

6.431x Probability – The Science of Uncertainty and Data

This is a cheat sheet for probability based on the online course given by Prof. John Tsitsiklis and Prof. Patrick Jaillet. Compiled by Janus B. Advincula.

Last Updated January 18, 2020

Probability Models and Axioms

Sample Space

We begin by listing the possible outcomes, Ω , of an experiment, i.e., the sample space. The list must be:

- mutually exclusive,
- collectively exhaustive, and
- at the *right* granularity

Probability Axioms

An event, A , is a subset of the sample space. Probability is assigned to events with the following axioms:

- Nonnegativity: $\mathbb{P}(A) \geq 0$
- Normalization: $\mathbb{P}(\Omega) = 1$
- (Finite) Additivity:
If $A \cap B = \emptyset$, then $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B)$.

Consequences of the Axioms

- $\mathbb{P}(A) \leq 1$
- $\mathbb{P}(\emptyset) = 0$
- $\mathbb{P}(A) + \mathbb{P}(A^c) = 1$
- If $A \subset B$, then $\mathbb{P}(A) \leq \mathbb{P}(B)$
- **Union bound:** $\mathbb{P}(A \cup B) \leq \mathbb{P}(A) + \mathbb{P}(B)$

Countable Additivity Axiom If A_1, A_2, A_3, \dots is an infinite *sequence* of *disjoint* events, then

$$\mathbb{P}(A_1 \cup A_2 \cup \dots) = \mathbb{P}(A_1) + \mathbb{P}(A_2) + \dots$$

$$\text{or } \mathbb{P}\left(\bigcup_i A_i\right) = \sum_i \mathbb{P}(A_i)$$

Discrete Uniform Law Assume that Ω is finite and consists of n equally likely elements. Also, assume that $A \subset \Omega$ consists of k elements. Then,

$$\mathbb{P}(A) = \frac{k}{n}.$$

Mathematical Background

Sets A set is a collection of distinct elements.

$$x \in S \cup T \Leftrightarrow x \in S \text{ or } x \in T$$

$$x \in S \cap T \Leftrightarrow x \in S \text{ and } x \in T$$

De Morgan's Laws

$$\left(\bigcup_n S_n\right)^c = \bigcap_n S_n^c$$

$$\left(\bigcap_n S_n\right)^c = \bigcup_n S_n^c$$

Geometric Series

$$S = \sum_{i=0}^{\infty} \alpha^i = 1 + \alpha + \alpha^2 + \dots = \frac{1}{1 - \alpha}$$

Bonferroni's Inequality

$$\mathbb{P}(A_1 \cap \dots \cap A_n) \geq \mathbb{P}(A_1) + \dots + \mathbb{P}(A_n) - (n - 1)$$

Conditioning and Independence

Conditioning and Bayes' Rule

Conditional Probability We denote the probability of A , given that B occurred by $\mathbb{P}(A|B)$ and this is defined as

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}, \quad \mathbb{P}(B) > 0.$$

Conditional probabilities share properties of ordinary probabilities.

The multiplication rule Given the definition of conditional probability, we have

$$\mathbb{P}(A \cap B) = \mathbb{P}(B) \mathbb{P}(A|B) = \mathbb{P}(A) \mathbb{P}(B|A).$$

In general, we have

$$\mathbb{P}(A_1 \cap A_2 \cap \dots \cap A_n) = \mathbb{P}(A_1) \prod_{i=2}^n \mathbb{P}(A_i | A_1 \cap \dots \cap A_{i-1})$$

Total probability theorem We partition the sample space into A_1, A_2, A_3, \dots . Then,

$$\begin{aligned} \mathbb{P}(B) &= \mathbb{P}(B \cap A_1) + \mathbb{P}(B \cap A_2) + \mathbb{P}(B \cap A_3) + \dots \\ &= \mathbb{P}(A_1) \mathbb{P}(B|A_1) + \mathbb{P}(A_2) \mathbb{P}(B|A_2) + \mathbb{P}(A_3) \mathbb{P}(B|A_3) + \dots \\ &= \sum_i \mathbb{P}(A_i) \mathbb{P}(B|A_i) \end{aligned}$$

Bayes' Rule

$$\mathbb{P}(A_i|B) = \frac{\mathbb{P}(A_i) \mathbb{P}(B|A_i)}{\sum_j \mathbb{P}(A_j) \mathbb{P}(B|A_j)}$$

Independence

Independence of Two Events A and B are independent if knowing whether A occurred gives no information about whether B occurred. More formally, A and B (which have nonzero probability) are independent if and only if one of the following equivalent statements holds:

$$\begin{aligned} \mathbb{P}(A \cap B) &= \mathbb{P}(A) \mathbb{P}(B) \\ \mathbb{P}(A|B) &= \mathbb{P}(A) \\ \mathbb{P}(B|A) &= \mathbb{P}(B) \end{aligned}$$

Independence of Event Complements If A and B are independent, then A and B^c are independent:

$$\mathbb{P}(A \cap B^c) = \mathbb{P}(A) \mathbb{P}(B^c)$$

Conditional Independence A and B are conditionally independent given C if

$$\mathbb{P}(A \cap B|C) = \mathbb{P}(A|C) \mathbb{P}(B|C).$$

Conditional independence does not imply independence, and independence does not imply conditional independence.

Counting

Basic Counting Principle

For a selection that can be done in r stages, with n_i choices at each stage i , the number of possible selections is

$$n_1 n_2 \dots n_r.$$

Permutations The number of ways of ordering n distinct elements is

$$n! = 1 \cdot 2 \cdot 3 \cdot \dots \cdot n.$$

Combinations Given a set of n elements, the number of ways of constructing an *ordered* sequence of k *distinct* elements is

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

Subsets The number of subsets of $\{1, \dots, n\}$ is

$$\sum_{k=0}^n \binom{n}{k} = \binom{n}{0} + \dots + \binom{n}{n} = 2^n.$$

Partitions Given an n -element set and nonnegative integers n_1, n_2, \dots, n_r , whose sum is equal to n , the number of ways of partitioning the set into r disjoint subsets, with the i^{th} subset containing exactly n_i elements, is equal to

$$\frac{n!}{n_1! n_2! \dots n_r!}.$$

Discrete Random Variables

Random Variables

A random variable associates a value to every possible outcome. It can take discrete or continuous values.

Probability Mass Function (PMF) Gives the probability that a *discrete* random variable takes on the value x .

$$p_X(x) = \mathbb{P}(X = x) = \mathbb{P}(\{\omega \in \Omega : X(\omega) = x\})$$

The PMF satisfies

$$p_X(x) \geq 0 \text{ and } \sum_x p_X(x) = 1.$$

Bernoulli random variable A Bernoulli random variable X with parameter $0 \leq p \leq 1$, $X \sim \text{Ber}(p)$, takes the following values:

$$X = \begin{cases} 1, & \text{w.p. } p, \\ 0, & \text{w.p. } 1 - p. \end{cases}$$

Discrete uniform random variable A discrete uniform random variable X between a and b , $X \sim \text{Uni}[a, b]$, takes any of the values $\{a, a + 1, \dots, b\}$ with probability $\frac{1}{b - a + 1}$.

Binomial random variable A binomial random variable X with parameter n and $p \in [0, 1]$, $X \sim \text{Bin}(n, p)$, takes values in the set $\{0, 1, \dots, n\}$ and has PMF

$$p_X(k) = \binom{n}{k} p^k (1 - p)^{n - k}, \quad \text{for } k = 0, 1, \dots, n.$$

Geometric random variable A geometric random variable X with parameter $p \in [0, 1]$, $X \sim \text{Geo}(p)$, takes values in the set $\{1, 2, \dots\}$ with probability

$$p_X(k) = (1 - p)^{k-1} p.$$

Expectation Value The expectation value of a discrete random variable is defined as

$$\mathbb{E}[X] = \sum_x x p_X(x).$$

Expectation of a Bernoulli random variable The expected value of a Bernoulli r.v. with parameter p is

$$\mathbb{E}[X] = p.$$

Expectation of a Discrete Uniform random variable The expected value of a discrete uniform r.v. is

$$\mathbb{E}[X] = \frac{b + a}{2}.$$

Properties of Expectations

- If $X \geq 0$, then $\mathbb{E}[X] \geq 0$.
- If $a \leq X \leq b$, then $a \leq \mathbb{E}[X] \leq b$.
- If c is a constant, $\mathbb{E}[c] = c$.
- Let X be a r.v. and let $Y = g(X)$. Then,

$$\mathbb{E}[Y] = \mathbb{E}[g(X)] = \sum_x g(x)p_X(x).$$

Linearity of Expectation

$$\mathbb{E}[aX + b] = a\mathbb{E}[X] + b$$

Variance and Standard Deviation

Definition of Variance Variance is a measure of the spread of a PMF. For a random variable with mean $\mu = \mathbb{E}[X]$, it is defined as

$$\text{Var}(X) = \mathbb{E}[(X - \mu)^2] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2.$$

Properties

$$\text{Var}(aX + b) = a^2 \text{Var}(X)$$

Standard Deviation

$$\sigma_X = \sqrt{\text{Var}(X)}$$

Variance of a Bernoulli random variable

$$\text{Var}(X) = p(1 - p)$$

Variance of a Discrete Uniform random variable

$$\text{Var}(X) = \frac{1}{12}(b - a)(b - a + 1)$$

Conditional PMF and Expectation

Total Expectation Theorem Given a random variable X and events A_1, \dots, A_n , we have

$$\mathbb{E}[X] = \mathbb{P}(A_1)\mathbb{E}[X|A_1] + \dots + \mathbb{P}(A_n)\mathbb{E}[X|A_n].$$

Continuous Random Variables

Continuous Random Variables (CRVs)

Definition A random variable is continuous if it can be described by a PDF, $f_X(x)$, such that

$$\mathbb{P}(a \leq X \leq a + \delta) \approx f_X(a) \cdot \delta.$$

Expectation Assuming that $\int_{-\infty}^{\infty} |x|f_X(x)dx < \infty$, the expected value of a random variable X is

$$\mathbb{E}[X] = \int_{-\infty}^{\infty} xf_X(x)dx.$$

Properties

- If $X \geq 0$, then $\mathbb{E}[X] \geq 0$.
- If $a \leq X \leq b$, then $a \leq \mathbb{E}[X] \leq b$.
-

$$\mathbb{E}[g(X)] = \int_{-\infty}^{\infty} g(x)f_X(x)dx$$

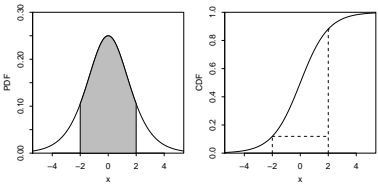
- $\mathbb{E}[aX + b] = a\mathbb{E}[X] + b$

What is the Probability Density Function (PDF)? The PDF f is the derivative of the CDF F .

$$F'(x) = f(x)$$

A PDF is nonnegative and integrates to 1. By the fundamental theorem of calculus, to get from PDF back to CDF we can integrate:

$$F(x) = \int_{-\infty}^x f(t)dt$$



To find the probability that a CRV takes on a value in an interval, integrate the PDF over that interval.

$$F(b) - F(a) = \int_a^b f(x)dx$$

Further Topics on Random Variables

Covariance and Correlation

Covariance is the analog of variance for two random variables.

$$\text{Cov}(X, Y) = E((X - E(X))(Y - E(Y))) = E(XY) - E(X)E(Y)$$

Note that

$$\text{Cov}(X, X) = E(X^2) - (E(X))^2 = \text{Var}(X)$$

Correlation is a standardized version of covariance that is always between -1 and 1 .

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$$

Covariance and Independence If two random variables are independent, then they are uncorrelated. The converse is not necessarily true (e.g., consider $X \sim \mathcal{N}(0, 1)$ and $Y = X^2$).

$$X \perp\!\!\!\perp Y \implies \text{Cov}(X, Y) = 0 \implies E(XY) = E(X)E(Y)$$

Covariance and Variance The variance of a sum can be found by

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$$

$$\text{Var}(X_1 + X_2 + \dots + X_n) = \sum_{i=1}^n \text{Var}(X_i) + 2 \sum_{i < j} \text{Cov}(X_i, X_j)$$

If X and Y are independent then they have covariance 0, so

$$X \perp\!\!\!\perp Y \implies \text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$$

If X_1, X_2, \dots, X_n are identically distributed and have the same covariance relationships (often by **symmetry**), then

$$\text{Var}(X_1 + X_2 + \dots + X_n) = n\text{Var}(X_1) + 2\binom{n}{2}\text{Cov}(X_1, X_2)$$

Covariance Properties For random variables W, X, Y, Z and constants a, b :

$$\begin{aligned} \text{Cov}(X, Y) &= \text{Cov}(Y, X) \\ \text{Cov}(X + a, Y + b) &= \text{Cov}(X, Y) \\ \text{Cov}(aX, bY) &= ab\text{Cov}(X, Y) \\ \text{Cov}(W + X, Y + Z) &= \text{Cov}(W, Y) + \text{Cov}(W, Z) + \text{Cov}(X, Y) \\ &\quad + \text{Cov}(X, Z) \end{aligned}$$

Correlation is location-invariant and scale-invariant For any constants a, b, c, d with a and c nonzero,

$$\text{Corr}(aX + b, cY + d) = \text{Corr}(X, Y)$$

Bayesian Inference

Consider an unknown random variable Θ with a prior distribution $p_{\Theta}(\theta)$ or $f_{\Theta}(\theta)$. We make an observation X modeled as a random variable with distribution $p_{X|\Theta}(x|\theta)$ or $f_{X|\Theta}(x|\theta)$. Using the appropriate version of the Bayes' rule, we can construct the appropriate posterior distribution $p_{\Theta|X}(\theta|x)$ or $f_{\Theta|X}(\theta|x)$.

Point Estimates

Maximum a posteriori probability (MAP) The MAP estimate, $\hat{\theta}$, is the value at which the posterior distribution is maximum:

$$p_{\Theta|X}(\theta^*|x) = \max_{\theta} p_{\Theta|X}(\theta|x)$$

$$f_{\Theta|X}(\theta^*|x) = \max_{\theta} f_{\Theta|X}(\theta|x).$$

Least Mean Squares (LMS) The LMS estimate is the conditional expectation of the posterior distribution:

$$\hat{\theta} = \mathbb{E}[\Theta|X = x].$$

Discrete Θ , Discrete X

$$p_{\Theta|X}(\theta|x) = \frac{p_{\Theta}(\theta)p_{X|\Theta}(x|\theta)}{p_X(x)}$$

$$p_X(x) = \sum_{\theta'} p_{\Theta}(\theta')p_{X|\Theta}(x|\theta')$$

Discrete Θ , Continuous X

$$p_{\Theta|X}(\theta|x) = \frac{p_{\Theta}(\theta)f_{X|\Theta}(x|\theta)}{f_X(x)}$$

$$f_X(x) = \sum_{\theta'} p_{\Theta}(\theta')f_{X|\Theta}(x|\theta')$$

Continuous Θ , Continuous X

$$f_{\Theta|X}(\theta|x) = \frac{f_{\Theta}(\theta)f_{X|\Theta}(x|\theta)}{f_X(x)}$$

$$f_X(x) = \int f_{\Theta}(\theta')f_{X|\Theta}(x|\theta')d\theta'$$

Continuous Θ , Discrete X

$$f_{\Theta|X}(\theta|x) = \frac{f_{\Theta}(\theta)p_{X|\Theta}(x|\theta)}{p_X(x)}$$

$$p_X(x) = \int f_{\Theta}(\theta')p_{X|\Theta}(x|\theta')d\theta'$$

Linear Least Mean Squares (LLMS) Estimation

In some cases, the conditional expectation $\mathbb{E}[\Theta|X]$ may be hard to compute or implement. In that case, we can restrict our attention to estimators of the form $\hat{\Theta} = aX + b$. Then,

$$\begin{aligned} \hat{\Theta}_{\text{LLMS}} &= \mathbb{E}[\Theta] + \frac{\text{Cov}(\Theta, X)}{\text{Var}(X)} (X - \mathbb{E}[X]) \\ &= \mathbb{E}[\Theta] + \rho \frac{\sigma_{\Theta}}{\sigma_X} (X - \mathbb{E}[X]) \end{aligned}$$

Limit Theorems and Classical Statistics

Markov Inequality

If $X \geq 0$ and $a > 0$, then

P(X ≥ a) ≤ E[X] / a.

Chebyshev Inequality

If the variance is small, then X is unlikely to be too far from the mean.

P(|X - μ| ≥ c) ≤ σ² / c²

The Weak Law of Large Numbers (WLLN)

Let X_1, X_2, X_3, \dots be i.i.d. with mean μ and variance σ^2 . The **sample mean** is

M_n = (X_1 + ... + X_n) / n.

The **Weak Law of Large Numbers** states that:

for ε > 0, P(|M_n - μ| ≥ ε) → 0, as n → ∞.

Transformations

One Variable Transformations Let's say that we have a random variable X with PDF $f_X(x)$, but we are also interested in some function of X . We call this function $Y = g(X)$. Also let $y = g(x)$. If g is differentiable and strictly increasing (or strictly decreasing), then the PDF of Y is

f_Y(y) = f_X(x) |dx/dy| = f_X(g⁻¹(y)) |d/dy g⁻¹(y)|

The derivative of the inverse transformation is called the **Jacobian**.

Two Variable Transformations Similarly, let's say we know the joint PDF of U and V but are also interested in the random vector (X, Y) defined by $(X, Y) = g(U, V)$. Let

∂(u, v) / ∂(x, y) = (∂u/∂x ∂u/∂y; ∂v/∂x ∂v/∂y)

be the **Jacobian matrix**. If the entries in this matrix exist and are continuous, and the determinant of the matrix is never 0, then

f_{X,Y}(x, y) = f_{U,V}(u, v) |∂(u, v) / ∂(x, y)|

The inner bars tells us to take the matrix's determinant, and the outer bars tell us to take the absolute value. In a 2 × 2 matrix,

| a b; c d | = |ad - bc|

Convolutions

Convolution Integral If you want to find the PDF of the sum of two independent CRVs X and Y , you can do the following integral:

f_{X+Y}(t) = ∫_{-∞}^∞ f_X(x) f_Y(t - x) dx

Example Let $X, Y \sim \mathcal{N}(0, 1)$ be i.i.d. Then for each fixed t ,

f_{X+Y}(t) = ∫_{-∞}^∞ 1/√(2π) e^{-x²/2} 1/√(2π) e^{-(t-x)²/2} dx

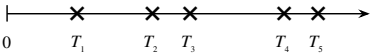
By completing the square and using the fact that a Normal PDF integrates to 1, this works out to $f_{X+Y}(t)$ being the $\mathcal{N}(0, 2)$ PDF.

Bernoulli and Poisson Processes

Definition We have a **Poisson process** of rate λ arrivals per unit time if the following conditions hold:

- 1. The number of arrivals in a time interval of length t is $\text{Pois}(\lambda t)$.
- 2. Numbers of arrivals in disjoint time intervals are independent.

For example, the numbers of arrivals in the time intervals $[0, 5]$, $(5, 12)$, and $[13, 23]$ are independent with $\text{Pois}(5\lambda)$, $\text{Pois}(7\lambda)$, $\text{Pois}(10\lambda)$ distributions, respectively.



Count-Time Duality Consider a Poisson process of emails arriving in an inbox at rate λ emails per hour. Let T_n be the time of arrival of the n th email (relative to some starting time 0) and N_t be the number of emails that arrive in $[0, t]$. Let's find the distribution of T_1 . The event $T_1 > t$, the event that you have to wait more than t hours to get the first email, is the same as the event $N_t = 0$, which is the event that there are no emails in the first t hours. So

P(T1 > t) = P(Nt = 0) = e^{-λt} → P(T1 ≤ t) = 1 - e^{-λt}

Thus we have $T_1 \sim \text{Expo}(\lambda)$. By the memoryless property and similar reasoning, the interarrival times between emails are i.i.d. $\text{Expo}(\lambda)$, i.e., the differences $T_n - T_{n-1}$ are i.i.d. $\text{Expo}(\lambda)$.

Order Statistics

Definition Let's say you have n i.i.d. r.v.s X_1, X_2, \dots, X_n . If you arrange them from smallest to largest, the i th element in that list is the i th order statistic, denoted $X_{(i)}$. So $X_{(1)}$ is the smallest in the list and $X_{(n)}$ is the largest in the list.

Note that the order statistics are *dependent*, e.g., learning $X_{(4)} = 42$ gives us the information that $X_{(1)}, X_{(2)}, X_{(3)}$ are ≤ 42 and $X_{(5)}, X_{(6)}, \dots, X_{(n)}$ are ≥ 42 .

Distribution Taking n i.i.d. random variables X_1, X_2, \dots, X_n with CDF $F(x)$ and PDF $f(x)$, the CDF and PDF of $X_{(i)}$ are:

F_{X_{(i)}}(x) = P(X_{(i)} ≤ x) = ∑_{k=i}^n (n choose k) F(x)^k (1 - F(x))^{n-k}

f_{X_{(i)}}(x) = n (n-1 choose i-1) F(x)^{i-1} (1 - F(x))^{n-i} f(x)

Uniform Order Statistics The j th order statistic of i.i.d. $U_1, \dots, U_n \sim \text{Unif}(0, 1)$ is $U_{(j)} \sim \text{Beta}(j, n - j + 1)$.

MVN, LLN, CLT

Central Limit Theorem (CLT)

Approximation using CLT

We use \sim to denote *is approximately distributed*. We can use the **Central Limit Theorem** to approximate the distribution of a random variable $Y = X_1 + X_2 + \dots + X_n$ that is a sum of n i.i.d. random variables X_i . Let $E(Y) = \mu_Y$ and $\text{Var}(Y) = \sigma_Y^2$. The CLT says

Y ~ N(μ_Y, σ_Y²)

If the X_i are i.i.d. with mean μ_X and variance σ_X^2 , then $\mu_Y = n\mu_X$ and $\sigma_Y^2 = n\sigma_X^2$. For the sample mean \bar{X}_n , the CLT says

̄X_n = 1/n (X_1 + X_2 + ... + X_n) ~ N(μ_X, σ_X²/n)

Asymptotic Distributions using CLT

We use \xrightarrow{D} to denote *converges in distribution to* as $n \rightarrow \infty$. The CLT says that if we standardize the sum $X_1 + \dots + X_n$ then the distribution of the sum converges to $\mathcal{N}(0, 1)$ as $n \rightarrow \infty$:

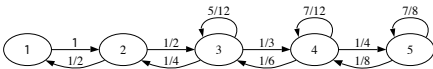
1 / (σ√n) (X_1 + ... + X_n - nμ_X) → N(0, 1)

In other words, the CDF of the left-hand side goes to the standard Normal CDF, Φ . In terms of the sample mean, the CLT says

√n(̄X_n - μ_X) / σ_X → N(0, 1)

Markov Chains

Definition



A Markov chain is a random walk in a **state space**, which we will assume is finite, say $\{1, 2, \dots, M\}$. We let X_t denote which element of the state space the walk is visiting at time t . The Markov chain is the sequence of random variables tracking where the walk is at all points in time, X_0, X_1, X_2, \dots . By definition, a Markov chain must satisfy the **Markov property**, which says that if you want to predict where the chain will be at a future time, if we know the present state then the entire past history is irrelevant. *Given the present, the past and future are conditionally independent*. In symbols,

P(X_{n+1} = j | X_0 = i_0, X_1 = i_1, ..., X_n = i) = P(X_{n+1} = j | X_n = i)

State Properties

A state is either recurrent or transient.

- If you start at a **recurrent state**, then you will always return back to that state at some point in the future. ♣ *You can check-out any time you like, but you can never leave.* ♣
- Otherwise you are at a **transient state**. There is some positive probability that once you leave you will never return. ♣ *You don't have to go home, but you can't stay here.* ♣

A state is either periodic or aperiodic.

- If you start at a **periodic state** of period k , then the GCD of the possible numbers of steps it would take to return back is $k > 1$.
- Otherwise you are at an **aperiodic state**. The GCD of the possible numbers of steps it would take to return back is 1.

Transition Matrix

Let the state space be $\{1, 2, \dots, M\}$. The transition matrix Q is the $M \times M$ matrix where element q_{ij} is the probability that the chain goes from state i to state j in one step:

q_{ij} = P(X_{n+1} = j | X_n = i)

To find the probability that the chain goes from state i to state j in exactly m steps, take the (i, j) element of Q^m .

q_{ij}^{(m)} = P(X_{n+m} = j | X_n = i)

If X_0 is distributed according to the row vector PMF \vec{p} , i.e., $p_j = P(X_0 = j)$, then the PMF of X_n is $\vec{p}Q^n$.

Chain Properties

A chain is **irreducible** if you can get from anywhere to anywhere. If a chain (on a finite state space) is irreducible, then all of its states are recurrent. A chain is **periodic** if any of its states are periodic, and is **aperiodic** if none of its states are periodic. In an irreducible chain, all states have the same period.

A chain is **reversible** with respect to \vec{s} if $s_i q_{ij} = s_j q_{ji}$ for all i, j . Examples of reversible chains include any chain with $q_{ij} = q_{ji}$, with $\vec{s} = (\frac{1}{M}, \frac{1}{M}, \dots, \frac{1}{M})$, and random walk on an undirected network.

Stationary Distribution

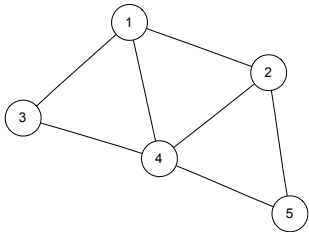
Let us say that the vector $\vec{s} = (s_1, s_2, \dots, s_M)$ be a PMF (written as a row vector). We will call \vec{s} the **stationary distribution** for the chain if $\vec{s}Q = \vec{s}$. As a consequence, if X_t has the stationary distribution, then all future X_{t+1}, X_{t+2}, \dots also have the stationary distribution.

For irreducible, aperiodic chains, the stationary distribution exists, is unique, and s_i is the long-run probability of a chain being at state i . The expected number of steps to return to i starting from i is $1/s_i$.

To find the stationary distribution, you can solve the matrix equation $(Q' - I)\vec{s}' = 0$. The stationary distribution is uniform if the columns of Q sum to 1.

Reversibility Condition Implies Stationarity If you have a PMF \vec{s} and a Markov chain with transition matrix Q , then $s_i q_{ij} = s_j q_{ji}$ for all states i, j implies that \vec{s} is stationary.

Random Walk on an Undirected Network



If you have a collection of **nodes**, pairs of which can be connected by undirected **edges**, and a Markov chain is run by going from the current node to a uniformly random node that is connected to it by an edge, then this is a random walk on an undirected network. The stationary distribution of this chain is proportional to the **degree sequence** (this is the sequence of degrees, where the degree of a node is how many edges are attached to it). For example, the stationary distribution of random walk on the network shown above is proportional to $(3, 3, 2, 4, 2)$, so it's $(\frac{3}{14}, \frac{3}{14}, \frac{2}{14}, \frac{4}{14}, \frac{2}{14})$.

Continuous Distributions

Uniform Distribution

Let us say that U is distributed $\text{Unif}(a, b)$. We know the following: **Properties of the Uniform** For a Uniform distribution, the probability of a draw from any interval within the support is proportional to the length of the interval. See *Universality of Uniform* and *Order Statistics* for other properties.

Example William throws darts really badly, so his darts are uniform over the whole room because they're equally likely to appear anywhere. William's darts have a Uniform distribution on the surface of the room. The Uniform is the only distribution where the probability of hitting in any specific region is proportional to the length/area/volume of that region,

and where the density of occurrence in any one specific spot is constant throughout the whole support.

Normal Distribution

Let us say that X is distributed $\mathcal{N}(\mu, \sigma^2)$. We know the following:

Central Limit Theorem The Normal distribution is ubiquitous because of the Central Limit Theorem, which states that the sample mean of i.i.d. r.v.s will approach a Normal distribution as the sample size grows, regardless of the initial distribution.

Location-Scale Transformation Every time we shift a Normal r.v. (by adding a constant) or rescale a Normal (by multiplying by a constant), we change it to another Normal r.v. For any Normal $X \sim \mathcal{N}(\mu, \sigma^2)$, we can transform it to the standard $\mathcal{N}(0, 1)$ by the following transformation:

Z = (X - μ) / σ ~ N(0, 1)

Standard Normal The Standard Normal, $Z \sim \mathcal{N}(0, 1)$, has mean 0 and variance 1. Its CDF is denoted by Φ .

Exponential Distribution

Let us say that X is distributed $\text{Expo}(\lambda)$. We know the following:

Story You're sitting on an open meadow right before the break of dawn, wishing that airplanes in the night sky were shooting stars, because you could really use a wish right now. You know that shooting stars come on average every 15 minutes, but a shooting star is not "due" to come just because you've waited so long. Your waiting time is memoryless; the additional time until the next shooting star comes does not depend on how long you've waited already.

Example The waiting time until the next shooting star is distributed $\text{Expo}(4)$ hours. Here $\lambda = 4$ is the **rate parameter**, since shooting stars arrive at a rate of 1 per 1/4 hour on average. The expected time until the next shooting star is $1/\lambda = 1/4$ hour.

Expos as a rescaled Expo(1)

Y ~ Expo(λ) → X = λY ~ Expo(1)

Memorylessness The Exponential Distribution is the only continuous memoryless distribution. The memoryless property says that for $X \sim \text{Expo}(\lambda)$ and any positive numbers s and t ,

P(X > s + t | X > s) = P(X > t)

Equivalently,

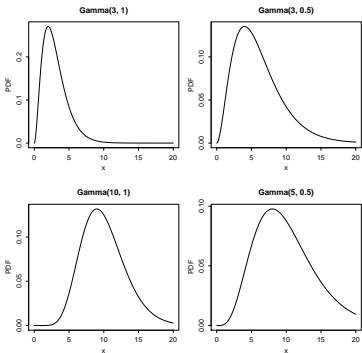
X - a | (X > a) ~ Expo(λ)

For example, a product with an $\text{Expo}(\lambda)$ lifetime is always "as good as new" (it doesn't experience wear and tear). Given that the product has survived a years, the additional time that it will last is still $\text{Expo}(\lambda)$.

Min of Expos If we have independent $X_i \sim \text{Expo}(\lambda_i)$, then $\min(X_1, \dots, X_k) \sim \text{Expo}(\lambda_1 + \lambda_2 + \dots + \lambda_k)$.

Max of Expos If we have i.i.d. $X_i \sim \text{Expo}(\lambda)$, then $\max(X_1, \dots, X_k)$ has the same distribution as $Y_1 + Y_2 + \dots + Y_k$, where $Y_j \sim \text{Expo}(j\lambda)$ and the Y_j are independent.

Gamma Distribution

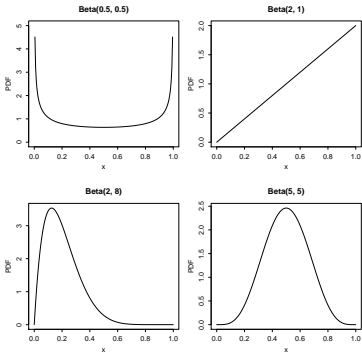


Let us say that X is distributed $\text{Gamma}(a, \lambda)$. We know the following:

Story You sit waiting for shooting stars, where the waiting time for a star is distributed $\text{Expo}(\lambda)$. You want to see n shooting stars before you go home. The total waiting time for the n th shooting star is $\text{Gamma}(n, \lambda)$.

Example You are at a bank, and there are 3 people ahead of you. The serving time for each person is Exponential with mean 2 minutes. Only one person at a time can be served. The distribution of your waiting time until it's your turn to be served is $\text{Gamma}(3, \frac{1}{2})$.

Beta Distribution



Conjugate Prior of the Binomial In the Bayesian approach to statistics, parameters are viewed as random variables, to reflect our uncertainty. The *prior* for a parameter is its distribution before observing data. The *posterior* is the distribution for the parameter after observing data. Beta is the *conjugate* prior of the Binomial because if you have a Beta-distributed prior on p in a Binomial, then the posterior distribution on p given the Binomial data is also Beta-distributed. Consider the following two-level model:

X|p ~ Bin(n, p)
p ~ Beta(a, b)

Then after observing $X = x$, we get the posterior distribution

p|(X = x) ~ Beta(a + x, b + n - x)

Order statistics of the Uniform See *Order Statistics*.

Beta-Gamma relationship If $X \sim \text{Gamma}(a, \lambda)$, $Y \sim \text{Gamma}(b, \lambda)$, with $X \perp\!\!\!\perp Y$ then

- $\frac{X}{X+Y} \sim \text{Beta}(a, b)$
- $X + Y \perp\!\!\!\perp \frac{X}{X+Y}$

This is known as the **bank–post office result**.

χ^2 (Chi-Square) Distribution

Let us say that X is distributed χ_n^2 . We know the following:

Story A Chi-Square(n) is the sum of the squares of n independent standard Normal r.v.s.

Properties and Representations

X is distributed as $Z_1^2 + Z_2^2 + \dots + Z_n^2$ for i.i.d. $Z_i \sim \mathcal{N}(0, 1)$

$X \sim \text{Gamma}(n/2, 1/2)$

Discrete Distributions

Distributions for four sampling schemes

	Replace	No Replace
Fixed # trials (n)	Binomial (Bern if $n = 1$)	HGeom
Draw until r success	NBin (Geom if $r = 1$)	NHGeom

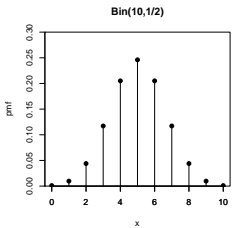
Bernoulli Distribution

The Bernoulli distribution is the simplest case of the Binomial distribution, where we only have one trial ($n = 1$). Let us say that X is distributed $\text{Bern}(p)$. We know the following:

Story A trial is performed with probability p of “success”, and X is the indicator of success: 1 means success, 0 means failure.

Example Let X be the indicator of Heads for a fair coin toss. Then $X \sim \text{Bern}(\frac{1}{2})$. Also, $1 - X \sim \text{Bern}(\frac{1}{2})$ is the indicator of Tails.

Binomial Distribution



Let us say that X is distributed $\text{Bin}(n, p)$. We know the following:

Story X is the number of “successes” that we will achieve in n independent trials, where each trial is either a success or a failure, each with the same probability p of success. We can also write X as a sum of multiple independent $\text{Bern}(p)$ random variables. Let $X \sim \text{Bin}(n, p)$ and $X_j \sim \text{Bern}(p)$, where all of the Bernoullis are independent. Then

$$X = X_1 + X_2 + X_3 + \dots + X_n$$

Example If Jeremy Lin makes 10 free throws and each one independently has a $\frac{3}{4}$ chance of getting in, then the number of free throws he makes is distributed $\text{Bin}(10, \frac{3}{4})$.

Properties Let $X \sim \text{Bin}(n, p)$, $Y \sim \text{Bin}(m, p)$ with $X \perp\!\!\!\perp Y$.

- Redefine success** $n - X \sim \text{Bin}(n, 1 - p)$
- Sum** $X + Y \sim \text{Bin}(n + m, p)$
- Conditional** $X|(X + Y = r) \sim \text{HGeom}(n, m, r)$
- Binomial-Poisson Relationship** $\text{Bin}(n, p)$ is approximately $\text{Pois}(\lambda)$ if p is small.
- Binomial-Normal Relationship** $\text{Bin}(n, p)$ is approximately $\mathcal{N}(np, np(1 - p))$ if n is large and p is not near 0 or 1.

Geometric Distribution

Let us say that X is distributed $\text{Geom}(p)$. We know the following:

Story X is the number of “failures” that we will achieve before we achieve our first success. Our successes have probability p .

Example If each pokeball we throw has probability $\frac{1}{10}$ to catch Mew, the number of failed pokeballs will be distributed $\text{Geom}(\frac{1}{10})$.

First Success Distribution

Equivalent to the Geometric distribution, except that it includes the first success in the count. This is 1 more than the number of failures. If $X \sim \text{FS}(p)$ then $E(X) = 1/p$.

Negative Binomial Distribution

Let us say that X is distributed $\text{NBin}(r, p)$. We know the following:

Story X is the number of “failures” that we will have before we achieve our r th success. Our successes have probability p .

Example Thundershock has 60% accuracy and can faint a wild Raticate in 3 hits. The number of misses before Pikachu faints Raticate with Thundershock is distributed $\text{NBin}(3, 0.6)$.

Hypergeometric Distribution

Let us say that X is distributed $\text{HGeom}(w, b, n)$. We know the following:

Story In a population of w desired objects and b undesired objects, X is the number of “successes” we will have in a draw of n objects, without replacement. The draw of n objects is assumed to be a **simple random sample** (all sets of n objects are equally likely).

Examples Here are some HGeom examples.

- Let’s say that we have only b Weedles (failure) and w Pikachus (success) in Viridian Forest. We encounter n Pokemon in the forest, and X is the number of Pikachus in our encounters.
- The number of Aces in a 5 card hand.
- You have w white balls and b black balls, and you draw n balls. You will draw X white balls.
- You have w white balls and b black balls, and you draw n balls without replacement. The number of white balls in your sample is $\text{HGeom}(w, b, n)$; the number of black balls is $\text{HGeom}(b, w, n)$.
- Capture-recapture** A forest has N elk, you capture n of them, tag them, and release them. Then you recapture a new sample of size m . How many tagged elk are now in the new sample? $\text{HGeom}(n, N - n, m)$

Poisson Distribution

Let us say that X is distributed $\text{Pois}(\lambda)$. We know the following:

Story There are rare events (low probability events) that occur many different ways (high possibilities of occurrences) at an average rate of λ occurrences per unit space or time. The number of events that occur in that unit of space or time is X .

Example A certain busy intersection has an average of 2 accidents per month. Since an accident is a low probability event that can happen many different ways, it is reasonable to model the number of accidents in a month at that intersection as $\text{Pois}(2)$. Then the number of accidents that happen in two months at that intersection is distributed $\text{Pois}(4)$.

Properties Let $X \sim \text{Pois}(\lambda_1)$ and $Y \sim \text{Pois}(\lambda_2)$, with $X \perp\!\!\!\perp Y$.

- Sum** $X + Y \sim \text{Pois}(\lambda_1 + \lambda_2)$
- Conditional** $X|(X + Y = n) \sim \text{Bin}\left(n, \frac{\lambda_1}{\lambda_1 + \lambda_2}\right)$
- Chicken-egg** If there are $Z \sim \text{Pois}(\lambda)$ items and we randomly and independently “accept” each item with probability p , then the number of accepted items $Z_1 \sim \text{Pois}(\lambda p)$, and the number of rejected items $Z_2 \sim \text{Pois}(\lambda(1 - p))$, and $Z_1 \perp\!\!\!\perp Z_2$.

Multivariate Distributions

Multinomial Distribution

Let us say that the vector $\vec{X} = (X_1, X_2, X_3, \dots, X_k) \sim \text{Mult}_k(n, \vec{p})$ where $\vec{p} = (p_1, p_2, \dots, p_k)$.

Story We have n items, which can fall into any one of the k buckets independently with the probabilities $\vec{p} = (p_1, p_2, \dots, p_k)$.

Example Let us assume that every year, 100 students in the Harry Potter Universe are randomly and independently sorted into one of four houses with equal probability. The number of people in each of the houses is distributed $\text{Mult}_4(100, \vec{p})$, where $\vec{p} = (0.25, 0.25, 0.25, 0.25)$. Note that $X_1 + X_2 + \dots + X_4 = 100$, and they are dependent.

Joint PMF For $n = n_1 + n_2 + \dots + n_k$,

$$P(\vec{X} = \vec{n}) = \frac{n!}{n_1!n_2! \dots n_k!} p_1^{n_1} p_2^{n_2} \dots p_k^{n_k}$$

Marginal PMF, Lumping, and Conditionals Marginally, $X_i \sim \text{Bin}(n, p_i)$ since we can define “success” to mean category i . If you lump together multiple categories in a Multinomial, then it is still Multinomial. For example, $X_i + X_j \sim \text{Bin}(n, p_i + p_j)$ for $i \neq j$ since we can define “success” to mean being in category i or j . Similarly, if $k = 6$ and we lump categories 1-2 and lump categories 3-5, then

$$(X_1 + X_2, X_3 + X_4 + X_5, X_6) \sim \text{Mult}_3(n, (p_1 + p_2, p_3 + p_4 + p_5, p_6))$$

Conditioning on some X_j also still gives a Multinomial:

$$X_1, \dots, X_{k-1} | X_k = n_k \sim \text{Mult}_{k-1}\left(n - n_k, \left(\frac{p_1}{1 - p_k}, \dots, \frac{p_{k-1}}{1 - p_k}\right)\right)$$

Variances and Covariances We have $X_i \sim \text{Bin}(n, p_i)$ marginally, so $\text{Var}(X_i) = np_i(1 - p_i)$. Also, $\text{Cov}(X_i, X_j) = -np_i p_j$ for $i \neq j$.

Multivariate Uniform Distribution

See the univariate Uniform for stories and examples. For the 2D Uniform on some region, probability is proportional to area. Every point in the support has equal density, of value $\frac{1}{\text{area of region}}$. For the 3D Uniform, probability is proportional to volume.

Multivariate Normal (MVN) Distribution

A vector $\vec{X} = (X_1, X_2, \dots, X_k)$ is Multivariate Normal if every linear combination is Normally distributed, i.e., $t_1X_1 + t_2X_2 + \dots + t_kX_k$ is Normal for any constants t_1, t_2, \dots, t_k . The parameters of the Multivariate Normal are the **mean vector** $\vec{\mu} = (\mu_1, \mu_2, \dots, \mu_k)$ and the **covariance matrix** where the (i, j) entry is $\text{Cov}(X_i, X_j)$.

Properties The Multivariate Normal has the following properties.

- Any subvector is also MVN.
- If any two elements within an MVN are uncorrelated, then they are independent.
- The joint PDF of a Bivariate Normal (X, Y) with $\mathcal{N}(0, 1)$ marginal distributions and correlation $\rho \in (-1, 1)$ is

f_{X,Y}(x,y) = \frac{1}{2\pi\tau} \exp\left(-\frac{1}{2\tau^2}(x^2 + y^2 - 2\rho xy)\right),

with \tau = \sqrt{1 - \rho^2}.

Distribution Properties

Important CDFs

Standard Normal Φ

Exponential(λ) $F(x) = 1 - e^{-\lambda x}$, for $x \in (0, \infty)$

Uniform(0,1) $F(x) = x$, for $x \in (0, 1)$

Convolutions of Random Variables

A convolution of n random variables is simply their sum. For the following results, let X and Y be *independent*.

- $X \sim \text{Pois}(\lambda_1), Y \sim \text{Pois}(\lambda_2) \longrightarrow X + Y \sim \text{Pois}(\lambda_1 + \lambda_2)$
- $X \sim \text{Bin}(n_1, p), Y \sim \text{Bin}(n_2, p) \longrightarrow X + Y \sim \text{Bin}(n_1 + n_2, p)$.
 $\text{Bin}(n, p)$ can be thought of as a sum of i.i.d. $\text{Bern}(p)$ r.v.s.
- $X \sim \text{Gamma}(a_1, \lambda), Y \sim \text{Gamma}(a_2, \lambda) \longrightarrow X + Y \sim \text{Gamma}(a_1 + a_2, \lambda)$. $\text{Gamma}(n, \lambda)$ with n an integer can be thought of as a sum of i.i.d. $\text{Expo}(\lambda)$ r.v.s.
- $X \sim \text{NBin}(r_1, p), Y \sim \text{NBin}(r_2, p) \longrightarrow X + Y \sim \text{NBin}(r_1 + r_2, p)$. $\text{NBin}(r, p)$ can be thought of as a sum of i.i.d. $\text{Geom}(p)$ r.v.s.
- $X \sim \mathcal{N}(\mu_1, \sigma_1^2), Y \sim \mathcal{N}(\mu_2, \sigma_2^2) \longrightarrow X + Y \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$

Special Cases of Distributions

- $\text{Bin}(1, p) \sim \text{Bern}(p)$
- $\text{Beta}(1, 1) \sim \text{Unif}(0, 1)$
- $\text{Gamma}(1, \lambda) \sim \text{Expo}(\lambda)$
- $\chi_n^2 \sim \text{Gamma}\left(\frac{n}{2}, \frac{1}{2}\right)$
- $\text{NBin}(1, p) \sim \text{Geom}(p)$

Inequalities

- Cauchy-Schwarz** $|E(XY)| \leq \sqrt{E(X^2)E(Y^2)}$
- Markov** $P(X \geq a) \leq \frac{E|X|}{a}$ for $a > 0$
- Chebyshev** $P(|X - \mu| \geq a) \leq \frac{\sigma^2}{a^2}$ for $E(X) = \mu, \text{Var}(X) = \sigma^2$
- Jensen** $E(g(X)) \geq g(E(X))$ for g convex; reverse if g is concave

Formulas

Geometric Series

1 + r + r^2 + \dots + r^{n-1} = \sum_{k=0}^{n-1} r^k = \frac{1 - r^n}{1 - r}

1 + r + r^2 + \dots = \frac{1}{1 - r} \text{ if } |r| < 1

Exponential Function (e^x)

e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!} = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots = \lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n

Gamma and Beta Integrals

You can sometimes solve complicated-looking integrals by pattern-matching to a gamma or beta integral:

\int_0^{\infty} x^{t-1} e^{-x} dx = \Gamma(t) \qquad \int_0^1 x^{a-1} (1-x)^{b-1} dx = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}

Also, $\Gamma(a + 1) = a\Gamma(a)$, and $\Gamma(n) = (n - 1)!$ if n is a positive integer.

Euler’s Approximation for Harmonic Sums

1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n} \approx \log n + 0.577 \dots

Stirling’s Approximation for Factorials

n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n

Miscellaneous Definitions

Medians and Quantiles Let X have CDF F . Then X has median m if $F(m) \geq 0.5$ and $P(X \geq m) \geq 0.5$. For X continuous, m satisfies $F(m) = 1/2$. In general, the a th quantile of X is $\min\{x : F(x) \geq a\}$; the median is the case $a = 1/2$.

log Statisticians generally use log to refer to natural log (i.e., base e).

i.i.d r.v.s Independent, identically-distributed random variables.

Distributions in R

Command	What it does
help(distributions)	shows documentation on distributions
dbinom(k,n,p)	PMF $P(X = k)$ for $X \sim \text{Bin}(n, p)$
pbinom(x,n,p)	CDF $P(X \leq x)$ for $X \sim \text{Bin}(n, p)$
qbinom(a,n,p)	a th quantile for $X \sim \text{Bin}(n, p)$
rbinom(r,n,p)	vector of r i.i.d. $\text{Bin}(n, p)$ r.v.s
dgeom(k,p)	PMF $P(X = k)$ for $X \sim \text{Geom}(p)$
dhypgeom(k,w,b,n)	PMF $P(X = k)$ for $X \sim \text{HGeom}(w, b, n)$
dnbinom(k,r,p)	PMF $P(X = k)$ for $X \sim \text{NBin}(r, p)$
dpois(k,r)	PMF $P(X = k)$ for $X \sim \text{Pois}(r)$
dbeta(x,a,b)	PDF $f(x)$ for $X \sim \text{Beta}(a, b)$
dchisq(x,n)	PDF $f(x)$ for $X \sim \chi_n^2$
dexp(x,b)	PDF $f(x)$ for $X \sim \text{Expo}(b)$
dgamma(x,a,r)	PDF $f(x)$ for $X \sim \text{Gamma}(a, r)$
dlnorm(x,m,s)	PDF $f(x)$ for $X \sim \mathcal{LN}(m, s^2)$
dnorm(x,m,s)	PDF $f(x)$ for $X \sim \mathcal{N}(m, s^2)$
dt(x,n)	PDF $f(x)$ for $X \sim t_n$
dunif(x,a,b)	PDF $f(x)$ for $X \sim \text{Unif}(a, b)$

The table above gives R commands for working with various named distributions. Commands analogous to `pbinom`, `qbinom`, and `rbinom` work for the other distributions in the table. For example, `pnorm`, `qnorm`, and `rnorm` can be used to get the CDF, quantiles, and random generation for the Normal. For the Multinomial, `dmultinom` can be used for calculating the joint PMF and `rmultinom` can be used for generating random vectors. For the Multivariate Normal, after installing and loading the `mvtnorm` package `dmvnorm` can be used for calculating the joint PDF and `rmvnorm` can be used for generating random vectors.

Recommended Resources

- Introduction to Probability Book (<http://bit.ly/introprobability>)
- Stat 110 Online (<http://stat110.net>)
- Stat 110 Quora Blog (<https://stat110.quora.com/>)
- Quora Probability FAQ (<http://bit.ly/probabilityfaq>)
- R Studio (<https://www.rstudio.com>)
- LaTeX File (github.com/wzchen/probability_cheatsheet)

Please share this cheatsheet with friends!

Table of Distributions

Distribution	PMF/PDF and Support	Expected Value	Variance	MGF
Bernoulli Bern(p)	$P(X = 1) = p$ $P(X = 0) = q = 1 - p$	p	pq	$q + pe^t$
Binomial Bin(n, p)	$P(X = k) = \binom{n}{k} p^k q^{n-k}$ $k \in \{0, 1, 2, \dots, n\}$	np	npq	$(q + pe^t)^n$
Geometric Geom(p)	$P(X = k) = q^k p$ $k \in \{0, 1, 2, \dots\}$	q/p	q/p^2	$\frac{p}{1-qe^t}, qe^t < 1$
Negative Binomial NBin(r, p)	$P(X = n) = \binom{r+n-1}{r-1} p^r q^n$ $n \in \{0, 1, 2, \dots\}$	rq/p	rq/p^2	$(\frac{p}{1-qe^t})^r, qe^t < 1$
Hypergeometric HGeom(w, b, n)	$P(X = k) = \binom{w}{k} \binom{b}{n-k} / \binom{w+b}{n}$ $k \in \{0, 1, 2, \dots, n\}$	$\mu = \frac{nw}{b+w}$	$\left(\frac{w+b-n}{w+b-1}\right) n \frac{\mu}{n} (1 - \frac{\mu}{n})$	messy
Poisson Pois(λ)	$P(X = k) = \frac{e^{-\lambda} \lambda^k}{k!}$ $k \in \{0, 1, 2, \dots\}$	λ	λ	$e^{\lambda(e^t-1)}$
Uniform Unif(a, b)	$f(x) = \frac{1}{b-a}$ $x \in (a, b)$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$	$\frac{e^{tb}-e^{ta}}{t(b-a)}$
Normal $\mathcal{N}(\mu, \sigma^2)$	$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/(2\sigma^2)}$ $x \in (-\infty, \infty)$	μ	σ^2	$e^{t\mu + \frac{\sigma^2 t^2}{2}}$
Exponential Expo(λ)	$f(x) = \lambda e^{-\lambda x}$ $x \in (0, \infty)$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$	$\frac{\lambda}{\lambda-t}, t < \lambda$
Gamma Gamma(a, λ)	$f(x) = \frac{1}{\Gamma(a)} (\lambda x)^a e^{-\lambda x} \frac{1}{x}$ $x \in (0, \infty)$	$\frac{a}{\lambda}$	$\frac{a}{\lambda^2}$	$\left(\frac{\lambda}{\lambda-t}\right)^a, t < \lambda$
Beta Beta(a, b)	$f(x) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$ $x \in (0, 1)$	$\mu = \frac{a}{a+b}$	$\frac{\mu(1-\mu)}{(a+b+1)}$	messy
Log-Normal $\mathcal{LN}(\mu, \sigma^2)$	$\frac{1}{x\sigma\sqrt{2\pi}} e^{-(\log x - \mu)^2/(2\sigma^2)}$ $x \in (0, \infty)$	$\theta = e^{\mu + \sigma^2/2}$	$\theta^2(e^{\sigma^2} - 1)$	doesn't exist
Chi-Square χ_n^2	$\frac{1}{2^{n/2}\Gamma(n/2)} x^{n/2-1} e^{-x/2}$ $x \in (0, \infty)$	n	$2n$	$(1-2t)^{-n/2}, t < 1/2$
Student- t t_n	$\frac{\Gamma((n+1)/2)}{\sqrt{n\pi}\Gamma(n/2)} (1+x^2/n)^{-(n+1)/2}$ $x \in (-\infty, \infty)$	0 if $n > 1$	$\frac{n}{n-2}$ if $n > 2$	doesn't exist