

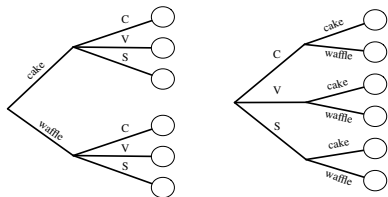
# Probability Cheatsheet

Compiled by William Chen (<http://wzchen.com>) and Joe Blitzstein, with contributions from Sebastian Chiu, Yuan Jiang, Yuqi Hou, and Jessy Hwang. Material based on Joe Blitzstein's (@stat110) lectures (<http://stat110.net>) and Blitzstein/Hwang's Introduction to Probability textbook (<http://bit.ly/introprobability>). Licensed under CC BY-NC-SA 4.0. Please share comments, suggestions, and errors at [http://github.com/wzchen/probability\\_cheatsheet](http://github.com/wzchen/probability_cheatsheet).

Last Updated August 22, 2020

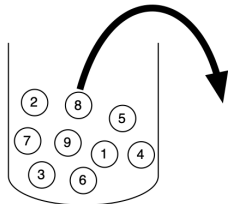
## Counting

### Multiplication Rule



Let's say we have a compound experiment (an experiment with multiple components). If the 1st component has  $n_1$  possible outcomes, the 2nd component has  $n_2$  possible outcomes, ..., and the  $r$ th component has  $n_r$  possible outcomes, then overall there are  $n_1 n_2 \dots n_r$  possibilities for the whole experiment.

### Sampling Table



The sampling table gives the number of possible samples of size  $k$  out of a population of size  $n$ , under various assumptions about how the sample is collected.

	Order Matters	Not Matter
With Replacement	$n^k$	$\binom{n+k-1}{k}$
Without Replacement	$\frac{n!}{(n-k)!}$	$\binom{n}{k}$

### Naive Definition of Probability

If all outcomes are equally likely, the probability of an event  $A$  happening is:

$$P_{\text{naive}}(A) = \frac{\text{number of outcomes favourable to } A}{\text{number of outcomes}}$$

## Thinking Conditionally

### Independence

**Independent 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:

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

**Conditional Independence**  $A$  and  $B$  are conditionally independent given  $C$  if  $P(A \cap B|C) = P(A|C)P(B|C)$ . Conditional independence does not imply independence, and independence does not imply conditional independence.

### Unions, Intersections, and Complements

**De Morgan's Laws** A useful identity that can make calculating probabilities of unions easier by relating them to intersections, and vice versa. Analogous results hold with more than two sets.

$$\begin{aligned}(A \cup B)^c &= A^c \cap B^c \\ (A \cap B)^c &= A^c \cup B^c\end{aligned}$$

### Joint, Marginal, and Conditional

**Joint Probability**  $P(A \cap B)$  or  $P(A, B)$  – Probability of  $A$  and  $B$ .

**Marginal (Unconditional) Probability**  $P(A)$  – Probability of  $A$ .

**Conditional Probability**  $P(A|B) = P(A, B)/P(B)$  – Probability of  $A$ , given that  $B$  occurred.

**Conditional Probability is Probability**  $P(A|B)$  is a probability function for any fixed  $B$ . Any theorem that holds for probability also holds for conditional probability.

### Probability of an Intersection or Union

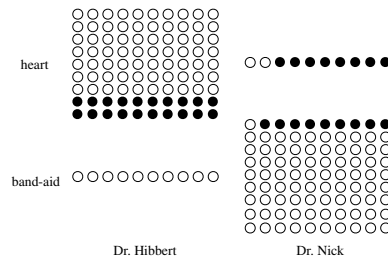
**Intersections via Conditioning**

$$\begin{aligned}P(A, B) &= P(A)P(B|A) \\ P(A, B, C) &= P(A)P(B|A)P(C|A, B)\end{aligned}$$

**Unions via Inclusion-Exclusion**

$$\begin{aligned}P(A \cup B) &= P(A) + P(B) - P(A \cap B) \\ P(A \cup B \cup C) &= P(A) + P(B) + P(C) \\ &\quad - P(A \cap B) - P(A \cap C) - P(B \cap C) \\ &\quad + P(A \cap B \cap C).\end{aligned}$$

### Simpson's Paradox



It is possible to have

$$\begin{aligned}P(A | B, C) &< P(A | B^c, C) \text{ and } P(A | B, C^c) < P(A | B^c, C^c) \\ \text{yet also } P(A | B) &> P(A | B^c).\end{aligned}$$

### Law of Total Probability (LOTP)

Let  $B_1, B_2, B_3, \dots, B_n$  be a *partition* of the sample space (i.e., they are disjoint and their union is the entire sample space).

$$\begin{aligned}P(A) &= P(A|B_1)P(B_1) + P(A|B_2)P(B_2) + \dots + P(A|B_n)P(B_n) \\ P(A) &= P(A \cap B_1) + P(A \cap B_2) + \dots + P(A \cap B_n)\end{aligned}$$

For **LOTP with extra conditioning**, just add in another event  $C$ !

$$\begin{aligned}P(A|C) &= P(A|B_1, C)P(B_1|C) + \dots + P(A|B_n, C)P(B_n|C) \\ P(A|C) &= P(A \cap B_1|C) + P(A \cap B_2|C) + \dots + P(A \cap B_n|C)\end{aligned}$$

Special case of LOTP with  $B$  and  $B^c$  as partition:

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

### Bayes' Rule

**Bayes' Rule, and with extra conditioning (just add in  $C$ !)**

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$P(A|B, C) = \frac{P(B|A, C)P(A|C)}{P(B|C)}$$

We can also write

$$P(A|B, C) = \frac{P(A, B, C)}{P(B, C)} = \frac{P(B, C|A)P(A)}{P(B, C)}$$

**Odds Form of Bayes' Rule**

$$\frac{P(A|B)}{P(A^c|B)} = \frac{P(B|A)}{P(B|A^c)} \frac{P(A)}{P(A^c)}$$

The *posterior odds* of  $A$  are the *likelihood ratio* times the *prior odds*.

## Random Variables and their Distributions

### PMF, CDF, and Independence

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

$$p_X(x) = P(X = x)$$

The PMF satisfies

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

**Cumulative Distribution Function (CDF)** Gives the probability that a random variable is less than or equal to  $x$ .

$$F_X(x) = P(X \leq x)$$

The CDF is an increasing, right-continuous function with

$$F_X(x) \rightarrow 0 \text{ as } x \rightarrow -\infty \text{ and } F_X(x) \rightarrow 1 \text{ as } x \rightarrow \infty$$

**Independence** Intuitively, two random variables are independent if knowing the value of one gives no information about the other. Discrete r.v.s  $X$  and  $Y$  are independent if for *all* values of  $x$  and  $y$

$$P(X = x, Y = y) = P(X = x)P(Y = y)$$

# Expected Value and Indicators

## Expected Value and Linearity

**Expected Value** (a.k.a. *mean*, *expectation*, or *average*) is a weighted average of the possible outcomes of our random variable. Mathematically, if  $x_1, x_2, x_3, \dots$  are all of the distinct possible values that  $X$  can take, the expected value of  $X$  is

$$E(X) = \sum_i x_i P(X = x_i)$$

$X$	$Y$	$X + Y$
3	4	7
2	2	4
6	8	14
10	23	33
1	-3	-2
1	0	1
5	9	14
4	1	5
...	...	...

$$\frac{1}{n} \sum_{i=1}^n x_i + \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n (x_i + y_i)$$
$$E(X) + E(Y) = E(X + Y)$$

**Linearity** For any r.v.s  $X$  and  $Y$ , and constants  $a, b, c$ ,

$$E(aX + bY + c) = aE(X) + bE(Y) + c$$

**Same distribution implies same mean** If  $X$  and  $Y$  have the same distribution, then  $E(X) = E(Y)$  and, more generally,

$$E(g(X)) = E(g(Y))$$

**Conditional Expected Value** is defined like expectation, only conditioned on any event  $A$ .

$$E(X|A) = \sum_x xP(X = x|A)$$

## Indicator Random Variables

**Indicator Random Variable** is a random variable that takes on the value 1 or 0. It is always an indicator of some event: if the event occurs, the indicator is 1; otherwise it is 0. They are useful for many problems about counting how many events of some kind occur. Write

$$I_A = \begin{cases} 1 & \text{if } A \text{ occurs,} \\ 0 & \text{if } A \text{ does not occur.} \end{cases}$$

Note that  $I_A^2 = I_A$ ,  $I_A I_B = I_{A \cap B}$ , and  $I_{A \cup B} = I_A + I_B - I_A I_B$ .

**Distribution**  $I_A \sim \text{Bern}(p)$  where  $p = P(A)$ .

**Fundamental Bridge** The expectation of the indicator for event  $A$  is the probability of event  $A$ :  $E(I_A) = P(A)$ .

## Variance and Standard Deviation

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

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

# Continuous RVs, LOTUS, UoU

## Continuous Random Variables (CRVs)

**What's the probability that a CRV is in an interval?** Take the difference in CDF values (or use the PDF as described later).

$$P(a \leq X \leq b) = P(X \leq b) - P(X \leq a) = F_X(b) - F_X(a)$$

For  $X \sim \mathcal{N}(\mu, \sigma^2)$ , this becomes

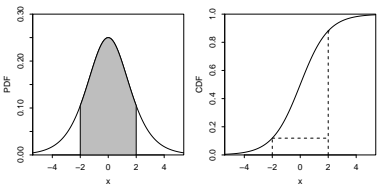
$$P(a \leq X \leq b) = \Phi\left(\frac{b - \mu}{\sigma}\right) - \Phi\left(\frac{a - \mu}{\sigma}\right)$$

**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.

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

**How do I find the expected value of a CRV?** Analogous to the discrete case, where you sum  $x$  times the PMF, for CRVs you integrate  $x$  times the PDF.

$$E(X) = \int_{-\infty}^{\infty} x f(x) dx$$

## LOTUS

**Expected value of a function of an r.v.** The expected value of  $X$  is defined this way:

$$E(X) = \sum_x xP(X = x) \text{ (for discrete } X)$$

$$E(X) = \int_{-\infty}^{\infty} x f(x) dx \text{ (for continuous } X)$$

The **Law of the Unconscious Statistician (LOTUS)** states that you can find the expected value of a *function of a random variable*,  $g(X)$ , in a similar way, by replacing the  $x$  in front of the PMF/PDF by  $g(x)$  but still working with the PMF/PDF of  $X$ :

$$E(g(X)) = \sum_x g(x)P(X = x) \text{ (for discrete } X)$$

$$E(g(X)) = \int_{-\infty}^{\infty} g(x)f(x)dx \text{ (for continuous } X)$$

**What's a function of a random variable?** A function of a random variable is also a random variable. For example, if  $X$  is the number of bikes you see in an hour, then  $g(X) = 2X$  is the number of bike wheels you see in that hour and  $h(X) = \binom{X}{2} = \frac{X(X-1)}{2}$  is the number of *pairs* of bikes such that you see both of those bikes in that hour.

**What's the point?** You don't need to know the PMF/PDF of  $g(X)$  to find its expected value. All you need is the PMF/PDF of  $X$ .

# Universality of Uniform (UoU)

When you plug any CRV into its own CDF, you get a Uniform(0,1) random variable. When you plug a Uniform(0,1) r.v. into an inverse CDF, you get an r.v. with that CDF. For example, let's say that a random variable  $X$  has CDF

$$F(x) = 1 - e^{-x}, \text{ for } x > 0$$

By UoU, if we plug  $X$  into this function then we get a uniformly distributed random variable.

$$F(X) = 1 - e^{-X} \sim \text{Unif}(0, 1)$$

Similarly, if  $U \sim \text{Unif}(0, 1)$  then  $F^{-1}(U)$  has CDF  $F$ . The key point is that for any continuous random variable  $X$ , we can transform it into a Uniform random variable and back by using its CDF.

# Moments and MGFs

## Moments

Moments describe the shape of a distribution. Let  $X$  have mean  $\mu$  and standard deviation  $\sigma$ , and  $Z = (X - \mu)/\sigma$  be the *standardized* version of  $X$ . The  $k$ th moment of  $X$  is  $\mu_k = E(X^k)$  and the  $k$ th standardized moment of  $X$  is  $m_k = E(Z^k)$ . The mean, variance, skewness, and kurtosis are important summaries of the shape of a distribution.

**Mean**  $E(X) = \mu_1$

**Variance**  $\text{Var}(X) = \mu_2 - \mu_1^2$

**Skewness**  $\text{Skew}(X) = m_3$

**Kurtosis**  $\text{Kurt}(X) = m_4 - 3$

## Moment Generating Functions

**MGF** For any random variable  $X$ , the function

$$M_X(t) = E(e^{tX})$$

is the **moment generating function (MGF)** of  $X$ , if it exists for all  $t$  in some open interval containing 0. The variable  $t$  could just as well have been called  $u$  or  $v$ . It's a bookkeeping device that lets us work with the *function*  $M_X$  rather than the *sequence* of moments.

**Why is it called the Moment Generating Function?** Because the  $k$ th derivative of the moment generating function, evaluated at 0, is the  $k$ th moment of  $X$ .

$$\mu_k = E(X^k) = M_X^{(k)}(0)$$

This is true by Taylor expansion of  $e^{tX}$  since

$$M_X(t) = E(e^{tX}) = \sum_{k=0}^{\infty} \frac{E(X^k)t^k}{k!} = \sum_{k=0}^{\infty} \frac{\mu_k t^k}{k!}$$

**MGF of linear functions** If we have  $Y = aX + b$ , then

$$M_Y(t) = E(e^{t(aX+b)}) = e^{bt} E(e^{(at)X}) = e^{bt} M_X(at)$$

**Uniqueness** *If it exists, the MGF uniquely determines the distribution.* This means that for any two random variables  $X$  and  $Y$ , they are distributed the same (their PMFs/PDFs are equal) if and only if their MGFs are equal.

**Summing Independent RVs by Multiplying MGFs.** If  $X$  and  $Y$  are independent, then

$$M_{X+Y}(t) = E(e^{t(X+Y)}) = E(e^{tX})E(e^{tY}) = M_X(t) \cdot M_Y(t)$$

The MGF of the sum of two random variables is the product of the MGFs of those two random variables.

# Joint PDFs and CDFs

## Joint Distributions

The **joint CDF** of  $X$  and  $Y$  is

$$F(x,y) = P(X \leq x, Y \leq y)$$

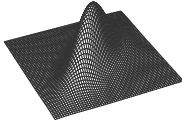
In the discrete case,  $X$  and  $Y$  have a **joint PMF**

$$p_{X,Y}(x,y) = P(X = x, Y = y).$$

In the continuous case, they have a **joint PDF**

$$f_{X,Y}(x,y) = \frac{\partial^2}{\partial x \partial y} F_{X,Y}(x,y).$$

The joint PMF/PDF must be nonnegative and sum/integrate to 1.



## Conditional Distributions

**Conditioning and Bayes’ rule for discrete r.v.s**

$$P(Y = y|X = x) = \frac{P(X = x, Y = y)}{P(X = x)} = \frac{P(X = x|Y = y)P(Y = y)}{P(X = x)}$$

**Conditioning and Bayes’ rule for continuous r.v.s**

$$f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{f_{X|Y}(x|y)f_Y(y)}{f_X(x)}$$

**Hybrid Bayes’ rule**

$$f_X(x|A) = \frac{P(A|X = x)f_X(x)}{P(A)}$$

## Marginal Distributions

To find the distribution of one (or more) random variables from a joint PMF/PDF, sum/integrate over the unwanted random variables.

**Marginal PMF from joint PMF**

$$P(X = x) = \sum_y P(X = x, Y = y)$$

**Marginal PDF from joint PDF**

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x,y)dy$$

## Independence of Random Variables

Random variables  $X$  and  $Y$  are independent if and only if any of the following conditions holds:

- Joint CDF is the product of the marginal CDFs
- Joint PMF/PDF is the product of the marginal PMFs/PDFs
- Conditional distribution of  $Y$  given  $X$  is the marginal distribution of  $Y$

Write  $X \perp\!\!\!\perp Y$  to denote that  $X$  and  $Y$  are independent.

## Multivariate LOTUS

LOTUS in more than one dimension is analogous to the 1D LOTUS. For discrete random variables:

$$E(g(X,Y)) = \sum_x \sum_y g(x,y)P(X = x, Y = y)$$

For continuous random variables:

$$E(g(X,Y)) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x,y)f_{X,Y}(x,y)dx dy$$

# Covariance and Transformations

## 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  $Y = X^2$ ).

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

**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,

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

## 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

$$f_Y(y) = f_X(x) \left| \frac{dx}{dy} \right| = f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right|$$

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

**Two Variable Transformations** Let  $X$  and  $Y$  be two jointly continuous random variables. Let

$$(Z,W) = g(X,Y) = (g_1(X,Y), g_2(X,Y))$$

where  $g : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  is a continuous one-to-one (invertible) function with continuous partial derivatives. Let  $h = g^{-1}$ , i.e.,

$$(X,Y) = h(Z,W) = (h_1(Z,W), h_2(Z,W))$$

Then  $Z$  and  $W$  are jointly continuous and their joint PDF,  $f_{ZW}(z,w)$ , for  $(z,w) \in R_{ZW}$  is given by

$$f_{ZW}(z,w) = f_{XY}(h_1(z,w), h_2(z,w))|J|$$

where  $J$  is the Jacobian of  $h$  defined by,

$$J = \det \begin{bmatrix} \frac{\partial h_1}{\partial z} & \frac{\partial h_1}{\partial w} \\ \frac{\partial h_2}{\partial z} & \frac{\partial h_2}{\partial w} \end{bmatrix} = \frac{\partial h_1}{\partial z} \cdot \frac{\partial h_2}{\partial w} - \frac{\partial h_2}{\partial z} \frac{\partial h_1}{\partial w}$$

If  $X \sim N(\mu_X, \sigma_X^2)$  and  $Y \sim N(\mu_Y, \sigma_Y^2)$  are independent, then they are jointly normal. Correlation matrix of  $X$ :  $\mathbf{R}_X = \mathbf{E}[\mathbf{X}\mathbf{X}^T]$  Covariance matrix of  $X$ :  $\mathbf{C}_X = \mathbf{E}[(\mathbf{X} - \mathbf{E}\mathbf{X})(\mathbf{X} - \mathbf{E}\mathbf{X})^T] = \mathbf{R}_X - \mathbf{E}\mathbf{X}\mathbf{E}\mathbf{X}^T$  Let  $G : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be a continuous and invertible function with continuous partial derivatives and let  $H = G^{-1}$ . Let  $\mathbf{Y} = G(\mathbf{X})$  and  $\mathbf{X} = G^{-1}(\mathbf{Y}) = H(\mathbf{Y}) = [H_1(\mathbf{Y}), \dots, H_n(\mathbf{Y})]^T$ .

## 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:

$$f_{X+Y}(t) = \int_{-\infty}^{\infty} f_X(x)f_Y(t-x)dx$$

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

$$f_{X+Y}(t) = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}}e^{-x^2/2} \frac{1}{\sqrt{2\pi}}e^{-(t-x)^2/2}dx$$

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.

# Random Vectors

Let  $\mathbf{X} = [X_1, X_2, \dots, X_n]^T$  and  $\mathbf{Y} = [Y_1, Y_2, \dots, Y_n]^T$  be random vectors. Then they have the following properties:

**CDF**  $F_{\mathbf{X}}(\mathbf{x}) = F_{X_1, \dots, X_n}(x_1, \dots, x_n) = P(X_1 \leq x_1, \dots, X_n \leq x_n)$

**PDF**  $f_{\mathbf{X}}(\mathbf{x}) = f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n)$  (for jointly continuous  $X_i$ ’s)

**Expectation**  $\mathbf{E}\mathbf{X} = [EX_1, EX_2, \dots, EX_n]^T$

**Correlation**  $\mathbf{R}_X = \mathbf{E}[\mathbf{X}\mathbf{X}^T]$  where  $[\mathbf{R}_X]_{ij} = E[X_i X_j]$ . Properties:

- $\mathbf{R}_X$  is **positive semidefinite**.
- $\mathbf{R}_X$  is **positive definite** if and only if all its eigenvalues are positive or equivalently  $\det(\mathbf{C}_X) > 0$ .

**Covariance**  $\mathbf{C}_X = \mathbf{E}[(\mathbf{X} - \mathbf{E}\mathbf{X})(\mathbf{X} - \mathbf{E}\mathbf{X})^T] = \mathbf{R}_X - \mathbf{E}\mathbf{X}\mathbf{E}\mathbf{X}^T$  where  $[\mathbf{C}_X]_{ij} = Cov[X_i, X_j] = E[X_i X_j] - E[X_i]E[X_j]$ .

**Cross-covariance**  $\mathbf{R}_{XY} = \mathbf{E}[\mathbf{X}\mathbf{Y}^T]$

**Cross-correlation**  $\mathbf{C}_{XY} = \mathbf{E}[(\mathbf{X} - \mathbf{E}\mathbf{X})(\mathbf{Y} - \mathbf{E}\mathbf{Y})^T]$

**Method of Transformations** Let  $\mathbf{Y} = G(\mathbf{X})$  where  $G : \mathbb{R}^n \mapsto \mathbb{R}^n$  is a continuous and invertible function with continuous partial derivatives and let  $H = G^{-1}$  so that  $\mathbf{X} = G^{-1}(\mathbf{Y}) = H(\mathbf{Y})$ . Then, the PDF of  $Y$  is given by

$$f_{\mathbf{Y}}(\mathbf{y}) = f_{\mathbf{X}}(H(\mathbf{y}))|J|$$

where  $J$  is the Jacobian of  $H$  defined by

$$J = \det \begin{bmatrix} \frac{\partial H_1}{\partial y_1} & \frac{\partial H_1}{\partial y_2} & \cdots & \frac{\partial H_1}{\partial y_n} \\ \frac{\partial H_2}{\partial y_1} & \frac{\partial H_2}{\partial y_2} & \cdots & \frac{\partial H_2}{\partial y_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial H_n}{\partial y_1} & \frac{\partial H_n}{\partial y_2} & \cdots & \frac{\partial H_n}{\partial y_n} \end{bmatrix}$$

**Linear Transformation** If  $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{b}$  then

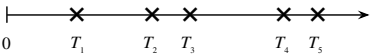
- 1.  $\mathbf{E}\mathbf{Y} = \mathbf{A}\mathbf{E}\mathbf{X} + \mathbf{b}$
- 2.  $\mathbf{C}_{\mathbf{Y}} = \mathbf{A}\mathbf{C}_{\mathbf{X}}\mathbf{A}^T$
- 3.  $f_{\mathbf{Y}}(\mathbf{y}) = \frac{1}{|\det \mathbf{A}|} f_{\mathbf{X}}(\mathbf{A}^{-1}(\mathbf{y} - \mathbf{b}))$

Poisson Process

**Definition** We have a **Poisson process** of rate  $\lambda$  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  $\lambda$  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

$$P(T_1 > t) = P(N_t = 0) = e^{-\lambda t} \implies P(T_1 \leq t) = 1 - e^{-\lambda t}$$

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  $\leq 42$  and  $X_{(5)}, X_{(6)}, \dots, X_{(n)}$  are  $\geq 42$ .

**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:

$$F_{X_{(i)}}(x) = P(X_{(i)} \leq x) = \sum_{k=i}^n \binom{n}{k} F(x)^k (1 - F(x))^{n-k}$$

$$f_{X_{(i)}}(x) = n \binom{n-1}{i-1} F(x)^{i-1} (1 - F(x))^{n-i} f(x)$$

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

Conditional Expectation

**Conditioning on an Event** We can find  $E(Y|A)$ , the expected value of  $Y$  given that event  $A$  occurred. A very important case is when  $A$  is the event  $X = x$ . Note that  $E(Y|A)$  is a *number*. For example:

- The expected value of a fair die roll, given that it is prime, is  $\frac{1}{3} \cdot 2 + \frac{1}{3} \cdot 3 + \frac{1}{3} \cdot 5 = \frac{10}{3}$ .
- Let  $Y$  be the number of successes in 10 independent Bernoulli trials with probability  $p$  of success. Let  $A$  be the event that the first 3 trials are all successes. Then

$$E(Y|A) = 3 + 7p$$

since the number of successes among the last 7 trials is  $\text{Bin}(7, p)$ .

- Let  $T \sim \text{Expo}(1/10)$  be how long you have to wait until the shuttle comes. Given that you have already waited  $t$  minutes, the expected additional waiting time is 10 more minutes, by the memoryless property. That is,  $E(T|T > t) = t + 10$ .

Discrete $Y$	Continuous $Y$
$E(Y) = \sum_y yP(Y = y)$	$E(Y) = \int_{-\infty}^{\infty} y f_Y(y) dy$
$E(Y A) = \sum_y yP(Y = y A)$	$E(Y A) = \int_{-\infty}^{\infty} y f(y A) dy$

**Conditioning on a Random Variable** We can also find  $E(Y|X)$ , the expected value of  $Y$  given the random variable  $X$ . This is a *function of the random variable  $X$* . It is *not* a number except in certain special cases such as if  $X \perp\!\!\!\perp Y$ . To find  $E(Y|X)$ , find  $E(Y|X = x)$  and then plug in  $X$  for  $x$ . For example:

- If  $E(Y|X = x) = x^3 + 5x$ , then  $E(Y|X) = X^3 + 5X$ .
- Let  $Y$  be the number of successes in 10 independent Bernoulli trials with probability  $p$  of success and  $X$  be the number of successes among the first 3 trials. Then  $E(Y|X) = X + 7p$ .
- Let  $X \sim \mathcal{N}(0, 1)$  and  $Y = X^2$ . Then  $E(Y|X = x) = x^2$  since if we know  $X = x$  then we know  $Y = x^2$ . And  $E(X|Y = y) = 0$  since if we know  $Y = y$  then we know  $X = \pm\sqrt{y}$ , with equal probabilities (by symmetry). So  $E(Y|X) = X^2, E(X|Y) = 0$ .

Properties of Conditional Expectation

- 1.  $E(Y|X) = E(Y)$  if  $X \perp\!\!\!\perp Y$
- 2.  $E(h(X)W|X) = h(X)E(W|X)$  (**taking out what's known**)  
In particular,  $E(h(X)|X) = h(X)$ .
- 3.  $E(E(Y|X)) = E(Y)$  (**Adam's Law**, a.k.a. Law of Total Expectation)

**Adam's Law (a.k.a. Law of Total Expectation)** can also be written in a way that looks analogous to LOTP. For any events  $A_1, A_2, \dots, A_n$  that partition the sample space,

$$E(Y) = E(Y|A_1)P(A_1) + \dots + E(Y|A_n)P(A_n)$$

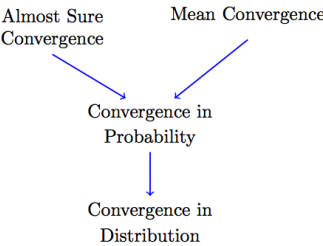
For the special case where the partition is  $A, A^c$ , this says

$$E(Y) = E(Y|A)P(A) + E(Y|A^c)P(A^c)$$

Eve's Law (a.k.a. Law of Total Variance)

$$\text{Var}(Y) = E(\text{Var}(Y|X)) + \text{Var}(E(Y|X))$$

Convergence of Random Variables



Convergence of a sequence of random variables  $X_1, X_2, X_3, \dots$  to random variable  $X$ .

Convergence in distribution

**Notation**  $X_n \xrightarrow{L^r} X$   
**Meaning**  $\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$ , for all  $x$  at which  $F_X(x)$  is continuous.

**Examples** *Central Limit Theorem*

Convergence in probability

**Notation**  $X_n \xrightarrow{L^r} X$  if,  
**Meaning**  $\lim_{n \rightarrow \infty} P(|X_n - X| \geq \epsilon) = 0$ , for all  $\epsilon > 0$

**Examples** *Weak Law of large Numbers*

Convergence in mean:

**Notation**  $X_n \xrightarrow{L^r} X$  if,  
**Meaning**  $\lim_{n \rightarrow \infty} E(|X_n - X|^r) = 0$

**mean-square convergence**  $X_n \xrightarrow{m.s.} X$  is defined as  $X_n \xrightarrow{L^2} X$

Almost sure convergence

**Notation**  $X_n \xrightarrow{a.s.} X$   
**Meaning**  $P(\{s \in S : \lim_{n \rightarrow \infty} X_n(s) = X(s)\}) = 1$ , where  $S$  is the sample space of  $X_i$  and  $X$

**Criteria** Criteria for a.s. convergence:

- $X_n \xrightarrow{a.s.} X$  if  $\sum_{n=1}^{\infty} P(|X_n - X| > \epsilon) < \infty$ , for all  $\epsilon > 0$
- $X_n \xrightarrow{a.s.} X$  iff  $\lim_{m \rightarrow \infty} P(A_m) = 1$ , for all  $\epsilon > 0$  where,  $A_m = \{|X_n - X| < \epsilon, \text{ for all } n \geq m\}$ .

**Examples** *Strong Law of large Numbers*

Continuous Mapping Theorem

- For a *continuous*  $h : \mathbb{R} \mapsto \mathbb{R}$
- $X_n \xrightarrow{d} X \implies h(X_n) \xrightarrow{d} h(X)$
- $X_n \xrightarrow{p} X \implies h(X_n) \xrightarrow{p} h(X)$
- $X_n \xrightarrow{a.s.} X \implies h(X_n) \xrightarrow{a.s.} h(X)$

Limit Theorems

Law of Large Numbers (LLN)

Let  $X_1, X_2, \dots, X_n$  be i.i.d. r.v.s and let be  $\bar{X}_n = \frac{X_1 + X_2 + \dots + X_n}{n}$  be the **sample mean**. Then the two variants of the **Law of Large Numbers** are:

**Weak (WLLN)**  $\bar{X}_n \xrightarrow{p} E[X_1]$

**Strong (SLLN)**  $\bar{X}_n \xrightarrow{a.s.} E[X_1]$

Central Limit Theorem (CLT)

Let  $X_1, X_2, \dots, X_n$  be i.i.d. with expected value  $E[X_i] = \mu < \infty$  and variance  $0 < Var[X_i] = \sigma^2 < \infty$  and let be  $\bar{X}_n = \frac{X_1 + X_2 + \dots + X_n}{n}$  be the **sample mean**. Then, the random variable,

Z\_n = (X\_bar - mu) / (sigma / sqrt(n)) = (X\_1 + X\_2 + ... + X\_n - n\*mu) / (sqrt(n)\*sigma)

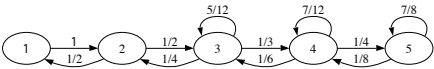
converges in distribution to the standard normal random variable as  $n \rightarrow \infty$ , i.e.  $Z_n \xrightarrow{d} \mathcal{N}(0, 1)$  or,

lim\_{n -> inf} P(Z\_n <= x) = Phi(x), for all x in R

where  $\Phi(x)$  is the standard normal CDF.

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^{(m)}\_{ij} = 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

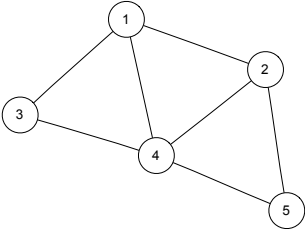
Let us say that the vector  $\vec{s} = (s_1, s_2, \dots, s_M)$  be a PMF (written as a row vector). We will call  $\vec{s}$  the **stationary distribution** for the chain if  $\vec{s}Q = \vec{s}$ . As a consequence, if  $X_t$  has the stationary distribution, then all future  $X_{t+1}, X_{t+2}, \dots$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

PDF f\_X(x) = 1 / (sigma \* sqrt(2\*pi)) \* e^(-(x - mu)^2 / (2\*sigma^2)) x in (-inf, inf)

Moments E[X] = mu Var[X] = sigma^2

**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 - mu) / sigma ~ 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  $\Phi$ .

Exponential Distribution

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

**Story** The amount of time until some specific event occurs, starting from now, being memoryless.

**Example** The waiting time until the next shooting star appears.

PDF f\_X(x) = lambda \* e^(-lambda \* x) x in (0, inf)

Moments E[X] = 1/lambda Var[X] = 1/lambda^2

**Expos as a rescaled Expo(1)**

Y ~ Expo(lambda) -> X = lambda \* Y ~ Expo(1)

**Memoryless** This 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,

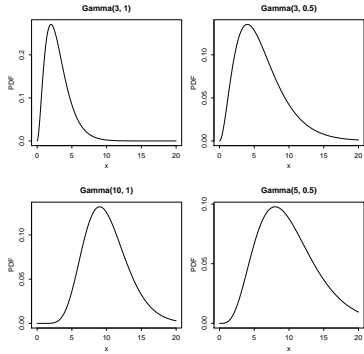
X - a | (X > a) ~ Expo(lambda)

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

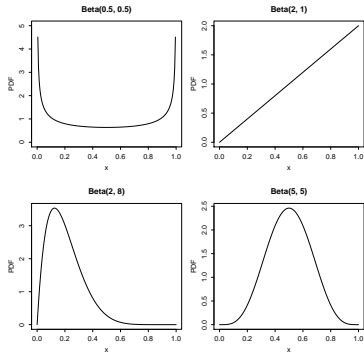


**Story** You sit waiting for shooting stars, where the waiting time for a star is distributed  $\text{Expo}(\lambda)$ . You want to see  $n$  shooting stars before you go home. The total waiting time for the  $n$ th shooting star is  $\text{Gamma}(n, \lambda)$ . If  $Z_1, Z_2, \dots, Z_n$  are independent r.v.s such that  $Z_i \sim \text{Expo}(\lambda)$ , then the r.v.  $X = Z_1 + Z_2 + \dots + Z_n$  is distributed  $\text{Gamma}(n, \lambda)$ .

**PDF**  $f_X(x) = \frac{1}{\Gamma(a)} (\lambda x)^a e^{-\lambda x} \frac{1}{x}, \quad x \in [0, \infty)$

**Moments**  $E[X] = \frac{a}{\lambda} \quad \text{Var}[X] = \frac{a}{\lambda^2}$

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:

$$X|p \sim \text{Bin}(n, p)$$
$$p \sim \text{Beta}(a, b)$$

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

$$p|(X = x) \sim \text{Beta}(a + x, b + n - x)$$

**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**.

Chi-Square ( $\chi^2$ ) Distribution

**Story** If  $Z_1, Z_2, \dots, Z_n$  are independent standard normal r.v.s, then the r.v.  $X = Z_1^2 + Z_2^2 + \dots + Z_n^2$  is said to have a chi-squared distribution with  $n$  *degrees of freedom* shown by  $X \sim \chi^2(n)$ . The chi-squared distribution is also a special case of the gamma distribution. More specifically,  $\chi^2(n) = \text{Gamma}(\frac{n}{2}, \frac{1}{2})$ .

**PDF**  $f_X(x) = \frac{\Gamma((n+1)/2)}{\sqrt{n\pi}\Gamma(n/2)} (1 + x^2/n)^{-(n+1)/2} \quad x \in (-\infty, \infty)$

**Moments**  $E[X] = n \quad \text{Var}[X] = 2n$

Student’s t-Distribution

**Story** If  $Z \sim \mathcal{N}(0, 1)$  and  $Y \sim \chi^2(n)$ , where  $n \in \mathbb{N}$  are independent r.v.s., then the r.v.  $X = Z/\sqrt{Y/n}$  is said to have a t-distribution with  $n$  *degrees of freedom* shown by  $X \sim T(n)$ .

**PDF**  $f_X(x) = \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} x^{\frac{n}{2}-1} e^{-\frac{x^2}{2}}, \quad x \in [0, \infty)$

**Moments**  $E[X] = 0$  if  $n > 1 \quad \text{Var}[X] = \frac{n}{n-2}$  if  $n > 2$

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:

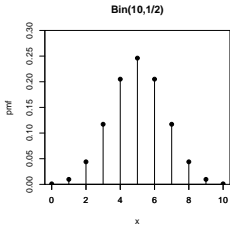
**Story** A trial is performed with probability  $p$  of “success”, and  $X$  is the indicator of success: 1 means success, 0 means failure.

**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.

**PMF**  $P(X = x) = p^x (1 - p)^{1-x}$

**Moments**  $E[X] = p \quad \text{Var}[X] = p(1 - p)$

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  $\text{Bin}(10, \frac{3}{4})$ .

**PMF**  $P_X(k) = \binom{n}{k} p^k (1 - p)^{n-k} \quad \text{for } k = 0, 1, 2, \dots, n$

**Moments**  $E[X] = np \quad \text{Var}[X] = np(1 - p)$

**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  $\text{Pois}(np)$  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 where the success probability is  $p$ .

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

Memoryless.

**PDF**  $P_X(k) = p(1 - p)^{k-1} \quad k \in 1, 2, 3, \dots$

**Moments**  $E[X] = \frac{1}{p} \quad \text{Var}[X] = \frac{1-p}{p^2}$

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  $\text{NBin}(3, 0.6)$ .

**PMF**  $P(X = k) = \binom{r+k-1}{k} p^r (1 - p)^k \quad k \in \{0, 1, 2, \dots\}$

**Moments**  $E[X] = \frac{r(1-p)}{p} \quad \text{Var}[X] = \frac{r(1-p)}{p^2}$

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).

**Example** 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  $\text{HGeom}(b, w, n)$ .

**PMF**  $P_X(k) = \binom{w}{k} \binom{b}{n-k} / \binom{w+b}{n} \quad k \in \{0, 1, 2, \dots, n\}$

**Moments**  $E[X] = \frac{nw}{w+b} = \mu \quad \text{Var}[X] = \left( \frac{w+b-n}{w+b-1} \right) n \frac{\mu}{n} (1 - \frac{\mu}{n})$

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  $\lambda$  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  $\text{Pois}(4)$ .

**PMF**  $P_X(k) = \frac{\lambda^k e^{-\lambda}}{k!} \quad k \in \{0, 1, 2, \dots\}$

**Moments**  $E[X] = \lambda \quad \text{Var}[X] = \lambda$

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

- 1. **Sum**  $X + Y \sim \text{Pois}(\lambda_1 + \lambda_2)$
- 2. **Conditional**  $X|(X + Y = n) \sim \text{Bin}\left(n, \frac{\lambda_1}{\lambda_1 + \lambda_2}\right)$
- 3. **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  $i \neq j$ .

Bivariate Normal Distribution

Two random variables  $X$  and  $Y$  are said to be *bivariate normal*, or *jointly normal*, if  $aX + bY$  has a normal distribution for all  $a, b \in \mathbb{R}$ . Jointly normal  $X$  and  $Y$  have the **Bivariate Normal distribution** with parameters  $\mu_X, \sigma_X^2, \mu_Y, \sigma_Y^2$ , and  $\rho$  and their joint PDF is given by

$$f_{XY}(x, y) = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \exp\left\{-\frac{1}{2(1-\rho^2)}\left[\left(\frac{x-\mu_X}{\sigma_X}\right)^2 + \left(\frac{y-\mu_Y}{\sigma_Y}\right)^2 - 2\rho\frac{(x-\mu_X)(y-\mu_Y)}{\sigma_X\sigma_Y}\right]\right\}$$

where  $\mu_X, \mu_Y \in \mathbb{R}, \sigma_X, \sigma_Y > 0$  and  $\rho \in (-1, 1)$  are the means, variances and correlation coefficient of  $X$  and  $Y$ . Properties of jointly normal  $X$  and  $Y$ :

- If  $X$  and  $Y$  are jointly normal random variables with parameters  $\mu_X, \sigma_X^2, \mu_Y, \sigma_Y^2$ , and  $\rho$ . Then, given  $X = x$ ,  $Y$  is normally distributed with

$$E[Y | X = x] = \mu_Y + \rho\sigma_Y \frac{x - \mu_X}{\sigma_X}$$
$$\text{Var}(Y | X = x) = (1 - \rho^2) \sigma_Y^2$$

- If  $X, Y \sim \mathcal{N}(0, 1)$ , then,

$$f_{XY}(x, y) = \frac{1}{2\pi\sqrt{1 - \rho^2}} \exp\left(-\frac{1}{2(1 - \rho^2)}(x^2 + y^2 - 2\rho xy)\right)$$

Multivariate Normal (MVN) Distribution

A random vector  $\mathbf{X} = (X_1, X_2, \dots, X_k)$  is said to be **Normal** or **Gaussian** if the random variables  $X_1, X_2, \dots, X_n$  are **jointly normal**

i.e. the random variable  $\sum_{i=1}^n a_i X_i$  is normal  $\forall a_i \in \mathbb{R}$ . A normal

random vector  $X$  has the **Multivariate Normal distribution**  $\mathcal{N}(\mathbf{m}, \mathbf{C})$  with mean  $\mathbf{m} = \mathbf{E}\mathbf{X}$  and covariance matrix  $\mathbf{C} = \mathbf{C}_\mathbf{x}$ . Its PDF is given by,

$$f_\mathbf{x}(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{n}{2}} \sqrt{\det \mathbf{C}}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1}(\mathbf{x} - \mathbf{m})\right\}$$

Further, a random vector  $\mathbf{Z} = (Z_1, Z_2, \dots, Z_k)$  is said to be **standard normal** if  $Z_i$ ’s are i.i.d. and  $Z_i \sim \mathcal{N}(0, 1)$ . A standard normal vector has  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  distribution (where  $\mathbf{I}$  is the identity matrix) i.e. its PDF is,

$$f_\mathbf{z}(\mathbf{z}) = \frac{1}{(2\pi)^{\frac{n}{2}}} \exp\left\{-\frac{1}{2}\mathbf{z}^T \mathbf{z}\right\}$$

Properties:

- Any subvector of a normal random vector is also normal.
- $\mathbf{X}$  is normal i.e.  $\mathbf{X} \sim \mathcal{N}(\mathbf{m}, \mathbf{C})$  iff there exists a standard normal  $\mathbf{Z}$  s.t.  $\mathbf{X} = \mathbf{AZ} + \mathbf{m}$  where  $\mathbf{A}$  is a matrix s.t.  $\mathbf{AA}^T = \mathbf{A}^T \mathbf{A} = \mathbf{C}$ . Note that if  $\text{svd}(\mathbf{C}) = \mathbf{QDQ}^T$  then  $\mathbf{A} = \mathbf{QD}^{\frac{1}{2}}\mathbf{Q}^T$ . This thus gives a way to generate arbitrary normal  $\mathbf{X}$  from a standard normal  $\mathbf{Z}$ .
- For a normal  $\mathbf{X} = [X_1, X_2, \dots, X_k]^T$ , following are equivalent:

- 1.  $X_i$ ’s are independent
- 2.  $X_i$ ’s are uncorrelated
- 3.  $\mathbf{C}_\mathbf{x}$  is diagonal.

Note that the above equivalence also holds for all subsets of  $X_i$ ’s as well.

- If  $\mathbf{X} \sim \mathcal{N}(\mathbf{m}, \mathbf{C})$  then  $\mathbf{Y} = \mathbf{AX} + \mathbf{b} \sim \mathcal{N}(\mathbf{AEX} + \mathbf{b}, \mathbf{ACA}^T)$

Distribution Properties

Important CDFs

**Standard Normal**  $\Phi$

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

**Uniform**( $\mathbf{0}, \mathbf{1}$ )  $F(x) = x$ , 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*.

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

Special Cases of Distributions

- 1.  $\text{Bin}(1, p) \sim \text{Bern}(p)$
- 2.  $\text{Beta}(1, 1) \sim \text{Unif}(0, 1)$
- 3.  $\text{Gamma}(1, \lambda) \sim \text{Expo}(\lambda)$
- 4.  $\chi^2(n) \sim \text{Gamma}(\frac{n}{2}, \frac{1}{2})$
- 5.  $\text{NBin}(1, p) \sim \text{Geom}(p)$

Inequalities

- 1. **Cauchy-Schwarz**  $|E(XY)| \leq \sqrt{E(X^2)E(Y^2)}$
- 2. **Union-bound**  $P\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n P(A_i)$  for any events  $A_1, A_2, \dots, A_n$ . Can be generalised using *Inclusion-Exclusion Principle*
- 3. **Markov**  $P(X \geq a) \leq \frac{E[X]}{a}$  for  $a > 0$
- 4. **Chebyshev**  $P(|X - E[X]| \geq a) \leq \frac{\text{Var}[X]}{a^2}$  for  $a > 0$
- 5. **Chernoff** Let  $M_X(s) = E[e^{sX}]$  be MGF. Then,  $P(X \geq a) \leq e^{-sa} M_X(s)$ , for all  $s > 0$   
 $P(X \leq a) \leq e^{-sa} M_X(s)$ , for all  $s < 0$
- 6. **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:

Γ(α) = ∫\_0^∞ t^{α-1} e^{-t} dt    B(α,β) = Γ(α)Γ(β) / Γ(α+β) = ∫\_0^1 t^{α-1} (1-t)^{β-1} dx

Also, Γ(α + 1) = αΓ(α), and Γ(n) = (n - 1)! if n is a positive integer.

Euler’s Approximation for Harmonic Sums

1 + 1/2 + 1/3 + ... + 1/n ≈ log n + 0.577 ...

Stirling’s Approximation for Factorials

n! ≈ √(2πn) (n/e)^n

Miscellaneous Definitions

Medians and Quantiles Let X have CDF F. Then X has median m if F(m) ≥ 0.5 and P(X ≥ m) ≥ 0.5. For X continuous, m satisfies F(m) = 1/2. In general, the ath quantile of X is min{x : F(x) ≥ 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.

Example Problems

Contributions from Sebastian Chiu

Calculating Probability

A textbook has n typos, which are randomly scattered amongst its n pages, independently. You pick a random page. What is the probability that it has no typos? Answer: There is a (1 - 1/n) probability that any specific typo isn’t on your page, and thus a

(1 - 1/n)^n probability that there are no typos on your page. For n

large, this is approximately e^{-1} = 1/e.

Linearity and Indicators (1)

In a group of n people, what is the expected number of distinct birthdays (month and day)? What is the expected number of birthday matches? Answer: Let X be the number of distinct birthdays and I\_j be the indicator for the jth day being represented.

E(I\_j) = 1 - P(no one born on day j) = 1 - (364/365)^n

By linearity, E(X) = 365 (1 - (364/365)^n). Now let Y be the number of birthday matches and J\_i be the indicator that the ith pair of people have the same birthday. The probability that any two specific people share a birthday is 1/365, so E(Y) = (n/2)/365.

Linearity and Indicators (2)

This problem is commonly known as the hat-matching problem. There are n people at a party, each with hat. At the end of the party, they each leave with a random hat. What is the expected number of people who leave with the right hat? Answer: Each hat has a 1/n chance of going to the right person. By linearity, the average number of hats that go to their owners is n(1/n) = 1.

Linearity and First Success

This problem is commonly known as the coupon collector problem. There are n coupon types. At each draw, you get a uniformly random coupon type. What is the expected number of coupons needed until you have a complete set? Answer: Let N be the number of coupons needed; we want E(N). Let N = N\_1 + ... + N\_n, where N\_1 is the draws to get our first new coupon, N\_2 is the additional draws needed to draw our second new coupon and so on. By the story of the First Success, N\_2 ~ FS((n - 1)/n) (after collecting first coupon type, there's (n - 1)/n chance you'll get something new). Similarly, N\_3 ~ FS((n - 2)/n), and N\_j ~ FS((n - j + 1)/n). By linearity,

E(N) = E(N\_1) + ... + E(N\_n) = n/n + n/(n-1) + ... + n/1 = n ∑\_{j=1}^n 1/j

This is approximately n(log(n) + 0.577) by Euler’s approximation.

Orderings of i.i.d. random variables

I call 2 UberX’s and 3 Lyfts at the same time. If the time it takes for the rides to reach me are i.i.d., what is the probability that all the Lyfts will arrive first? Answer: Since the arrival times of the five cars are i.i.d., all 5! orderings of the arrivals are equally likely. There are 3!2! orderings that involve the Lyfts arriving first, so the probability

that the Lyfts arrive first is 3!2!/5! = 1/10. Alternatively, there are 5/3 ways to choose 3 of the 5 slots for the Lyfts to occupy, where each of the choices are equally likely. One of these choices has all 3 of the

Lyfts arriving first, so the probability is 1/(5/3) = 1/10.

Expectation of Negative Hypergeometric

What is the expected number of cards that you draw before you pick your first Ace in a shuffled deck (not counting the Ace)? Answer: Consider a non-Ace. Denote this to be card j. Let I\_j be the indicator that card j will be drawn before the first Ace. Note that I\_j = 1 says that j is before all 4 of the Aces in the deck. The probability that this occurs is 1/5 by symmetry. Let X be the number of cards drawn before the first Ace. Then X = I\_1 + I\_2 + ... + I\_48, where each indicator corresponds to one of the 48 non-Aces. Thus,

E(X) = E(I\_1) + E(I\_2) + ... + E(I\_48) = 48/5 = 9.6.

Minimum and Maximum of RVs

What is the CDF of the maximum of n independent Unif(0,1) random variables? Answer: Note that for r.v.s X\_1, X\_2, ..., X\_n,

P(min(X\_1, X\_2, ..., X\_n) ≥ a) = P(X\_1 ≥ a, X\_2 ≥ a, ..., X\_n ≥ a)

Similarly,

P(max(X\_1, X\_2, ..., X\_n) ≤ a) = P(X\_1 ≤ a, X\_2 ≤ a, ..., X\_n ≤ a)

We will use this principle to find the CDF of U\_(n), where U\_(n) = max(U\_1, U\_2, ..., U\_n) and U\_i ~ Unif(0, 1) are i.i.d.

P(max(U\_1, U\_2, ..., U\_n) ≤ a) = P(U\_1 ≤ a, U\_2 ≤ a, ..., U\_n ≤ a) = P(U\_1 ≤ a)P(U\_2 ≤ a) ... P(U\_n ≤ a) = a^n

for 0 < a < 1 (and the CDF is 0 for a ≤ 0 and 1 for a ≥ 1).

Pattern-matching with e^x Taylor series

For X ~ Pois(λ), find E(1/(X + 1)). Answer: By LOTUS,

E(1/(X + 1)) = ∑\_{k=0}^∞ 1/(k + 1) \* e^{-λ} λ^k / k! = e^{-λ} / λ ∑\_{k=0}^∞ λ^{k+1} / (k + 1)! = (e^{-λ} / λ) (e^λ - 1)

Adam’s Law and Eve’s Law

William really likes speedsolving Rubik’s Cubes. But he’s pretty bad at it, so sometimes he fails. On any given day, William will attempt N ~ Geom(s) Rubik’s Cubes. Suppose each time, he has probability p of solving the cube, independently. Let T be the number of Rubik’s Cubes he solves during a day. Find the mean and variance of T.

Answer: Note that T|N ~ Bin(N, p). So by Adam’s Law,

E(T) = E(E(T|N)) = E(Np) = p(1 - s) / s

Similarly, by Eve’s Law, we have that

Var(T) = E(Var(T|N)) + Var(E(T|N)) = E(Np(1 - p)) + Var(Np) = p(1 - p)(1 - s) / s + p^2(1 - s) / s^2 = p(1 - s)(p + s(1 - p)) / s^2

MGF – Finding Moments

Find E(X^3) for X ~ Expo(λ) using the MGF of X. Answer: The MGF of an Expo(λ) is M(t) = λ / (λ - t). To get the third moment, we can take the third derivative of the MGF and evaluate at t = 0:

E(X^3) = 6 / λ^3

But a much nicer way to use the MGF here is via pattern recognition: note that M(t) looks like it came from a geometric series:

1 / (1 - t/λ) = ∑\_{n=0}^∞ (t/λ)^n = ∑\_{n=0}^∞ t^n / λ^n n!

The coefficient of t^n / n! here is the nth moment of X, so we have E(X^n) = n! / λ^n for all nonnegative integers n.

Markov chains (1)

Suppose X\_n is a two-state Markov chain with transition matrix

Q = [ [0, 1], [1 - α, α - β] ]

Find the stationary distribution s = (s\_0, s\_1) of X\_n by solving sQ = s, and show that the chain is reversible with respect to s. Answer: The equation sQ = s says that

s\_0 = s\_0(1 - α) + s\_1β and s\_1 = s\_0(α) + s\_0(1 - β)

By solving this system of linear equations, we have

s = ( β / (α + β), α / (α + β) )

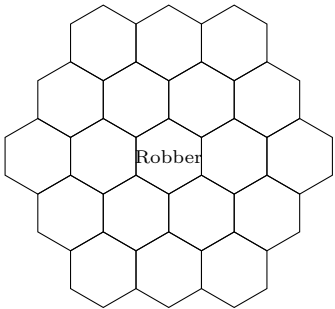
To show that the chain is reversible with respect to s, we must show s\_i q\_{ij} = s\_j q\_{ji} for all i, j. This is done if we can show s\_0 q\_{01} = s\_1 q\_{10}. And indeed,

s\_0 q\_{01} = αβ / (α + β) = s\_1 q\_{10}



Markov chains (2)

William and Sebastian play a modified game of Settlers of Catan, where every turn they randomly move the robber (which starts on the center tile) to one of the adjacent hexagons.



- (a) Is this Markov chain irreducible? Is it aperiodic? **Answer:** Yes to both. The Markov chain is irreducible because it can get from anywhere to anywhere else. The Markov chain is aperiodic because the robber can return back to a square in 2, 3, 4, 5, . . . moves, and the GCD of those numbers is 1.
- (b) What is the stationary distribution of this Markov chain? **Answer:** Since this is a random walk on an undirected graph, the stationary distribution is proportional to the degree sequence. The degree for the corner pieces is 3, the degree for the edge pieces is 4, and the degree for the center pieces is 6. To normalize this degree sequence, we divide by its sum. The sum of the degrees is  $6(3) + 6(4) + 7(6) = 84$ . Thus the stationary probability of being on a corner is  $3/84 = 1/28$ , on an edge is  $4/84 = 1/21$ , and in the center is  $6/84 = 1/14$ .
- (c) What fraction of the time will the robber be in the center tile in this game, in the long run? **Answer:** By the above,  $1/14$ .
- (d) What is the expected amount of moves it will take for the robber to return to the center tile? **Answer:** Since this chain is irreducible and aperiodic, to get the expected time to return we can just invert the stationary probability. Thus on average it will take  $14$  turns for the robber to return to the center tile.

Problem-Solving Strategies

Contributions from Jessy Hwang, Yuan Jiang, Yuqi Hou

1. **Getting started.** Start by *defining relevant events and random variables*. (“Let  $A$  be the event that I pick the fair coin”; “Let  $X$  be the number of successes.”) Clear notion is important for clear thinking! Then decide what it is that you’re supposed to be finding, in terms of your notation (“I want to find  $P(X = 3|A)$ ”). Think about what type of object your answer should be (a number? A random variable? A PMF? A PDF?) and what it should be in terms of. *Try simple and extreme cases*. To make an abstract experiment more concrete, try *drawing a picture* or making up numbers that could have happened. Pattern recognition: does the structure of the problem resemble something we’ve seen before?
2. **Calculating probability of an event.** Use counting principles if the naive definition of probability applies. Is the probability of the complement easier to find? Look for symmetries. Look for something to condition on, then apply Bayes’ Rule or the Law of Total Probability.

3. **Finding the distribution of a random variable.** First make sure you need the full distribution not just the mean (see next item). Check the *support* of the random variable: what values can it take on? Use this to rule out distributions that don’t fit. Is there a *story* for one of the named distributions that fits the problem at hand? Can you write the random variable as a function of an r.v. with a known distribution, say  $Y = g(X)$ ?
4. **Calculating expectation.** If it has a named distribution, check out the table of distributions. If it’s a function of an r.v. with a named distribution, try LOTUS. If it’s a count of something, try breaking it up into indicator r.v.s. If you can condition on something natural, consider using Adam’s law.
5. **Calculating variance.** Consider independence, named distributions, and LOTUS. If it’s a count of something, break it up into a sum of indicator r.v.s. If it’s a sum, use properties of covariance. If you can condition on something natural, consider using Eve’s Law.
6. **Calculating  $E(X^2)$ .** Do you already know  $E(X)$  or  $\text{Var}(X)$ ? Recall that  $\text{Var}(X) = E(X^2) - (E(X))^2$ . Otherwise try LOTUS.
7. **Calculating covariance.** Use the properties of covariance. If you’re trying to find the covariance between two components of a Multinomial distribution,  $X_i, X_j$ , then the covariance is  $-np_i p_j$  for  $i \neq j$ .
8. **Symmetry.** If  $X_1, \dots, X_n$  are i.i.d., consider using symmetry.
9. **Calculating probabilities of orderings.** Remember that all  $n!$  ordering of i.i.d. continuous random variables  $X_1, \dots, X_n$  are equally likely.
10. **Determining independence.** There are several equivalent definitions. Think about simple and extreme cases to see if you can find a counterexample.
11. **Do a painful integral.** If your integral looks painful, see if you can write your integral in terms of a known PDF (like Gamma or Beta), and use the fact that PDFs integrate to 1?
12. **Before moving on.** Check some simple and extreme cases, check whether the answer seems plausible, check for biohazards.

Biohazards

Contributions from Jessy Hwang

1. **Don’t misuse the naive definition of probability.** When answering “What is the probability that in a group of 3 people, no two have the same birth month?”, it is *not* correct to treat the people as indistinguishable balls being placed into 12 boxes, since that assumes the list of birth months {January, January, January} is just as likely as the list {January, April, June}, even though the latter is six times more likely.
2. **Don’t confuse unconditional, conditional, and joint probabilities.** In applying  $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$ , it is *not* correct to say “ $P(B) = 1$  because we know  $B$  happened”;  $P(B)$  is the *prior* probability of  $B$ . Don’t confuse  $P(A|B)$  with  $P(A, B)$ .
3. **Don’t assume independence without justification.** In the matching problem, the probability that card 1 is a match and card 2 is a match is not  $1/n^2$ . Binomial and Hypergeometric are often confused; the trials are independent in the Binomial story and dependent in the Hypergeometric story.
4. **Don’t forget to do sanity checks.** Probabilities must be between 0 and 1. Variances must be  $\geq 0$ . Supports must make sense. PMFs must sum to 1. PDFs must integrate to 1.

5. **Don’t confuse random variables, numbers, and events.** Let  $X$  be an r.v. Then  $g(X)$  is an r.v. for any function  $g$ . In particular,  $X^2, |X|, F(X)$ , and  $I_{X>3}$  are r.v.s.  $P(X^2 < X|X \geq 0), E(X), \text{Var}(X)$ , and  $g(E(X))$  are numbers.  $X = 2$  and  $F(X) \geq -1$  are events. It does not make sense to write  $\int_{-\infty}^{\infty} F(X)dx$ , because  $F(X)$  is a random variable. It does not make sense to write  $P(X)$ , because  $X$  is not an event.
6. **Don’t confuse a random variable with its distribution.** To get the PDF of  $X^2$ , you can’t just square the PDF of  $X$ . The right way is to use transformations. To get the PDF of  $X + Y$ , you can’t just add the PDF of  $X$  and the PDF of  $Y$ . The right way is to compute the convolution.
7. **Don’t pull non-linear functions out of expectations.**  $E(g(X))$  does not equal  $g(E(X))$  in general. The St. Petersburg paradox is an extreme example. See also Jensen’s inequality. The right way to find  $E(g(X))$  is with LOTUS.

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)	ath 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 $X \sim \text{Geom}(p)$
dhyp(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 $X \sim \chi^2(n)$
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](https://github.com/wzchen/probability_cheatsheet))

Please share this cheatsheet with friends!  
<http://wzchen.com/probability-cheatsheet>

# Table of Distributions

Distribution				
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( $r, p$ )	$P(X = k) = \binom{r+k-1}{n} p^r q^k$ $k \in \{0, 1, 2, \dots\}$	$rq/p$	$rq/p^2$	$(\frac{p}{1-qe^t})^r, qe^t < 1$
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( $\lambda$ )	$P(X = k) = \frac{\lambda^k e^{-\lambda}}{k!}$ $k \in \{0, 1, 2, \dots\}$	$\lambda$	$\lambda$	$e^{\lambda(e^t - 1)}$
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 $\mathcal{N}(\mu, \sigma^2)$	$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/(2\sigma^2)}$ $x \in (-\infty, \infty)$	$\mu$	$\sigma^2$	$e^{t\mu + \frac{\sigma^2 t^2}{2}}$
Exponential Expo( $\lambda$ )	$f(x) = \lambda e^{-\lambda x}$ $x \in (0, \infty)$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$	$\frac{\lambda}{\lambda - t}, t < \lambda$
Gamma Gamma( $a, \lambda$ )	$f(x) = \frac{1}{\Gamma(a)} (\lambda x)^a e^{-\lambda x} \frac{1}{x}$ $x \in (0, \infty)$	$\frac{a}{\lambda}$	$\frac{a}{\lambda^2}$	$\left(\frac{\lambda}{\lambda - t}\right)^a, t < \lambda$
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 $\chi^2(n)$	$\frac{1}{2^{n/2}\Gamma(n/2)} x^{n/2-1} e^{-x/2}$ $x \in (0, \infty)$	$n$	$2n$	$(1 - 2t)^{-n/2}, t < 1/2$
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