

# Random Variables

Fall 2020

Instructor:

Ajit Rajwade

# Topic Overview

- Random variable: definition
- Discrete and continuous random variables
- Probability density function (pdf) and cumulative distribution function (cdf)
- Joint and conditional pdfs
- Expectation and its properties
- Variance and covariance
- Markov's and Chebyshev's inequality
- Weak law of large numbers
- Moment generating functions

# Random variable

- In many random experiments, we are not always interested in the observed values, but in some numerical quantity determined by the observed values.
- Example: we may be interested in the sum of the values of two dice throws, or the number of heads appearing in  $n$  consecutive coin tosses.
- Any such quantities determined by the results of random experiments are called **random variables** (they may also be the observations themselves).

# Random variable

Value of $X$ (Denoted as $x$ ) where $X$ = sum of 2 dice throws	$P(X=x)$
2	1/36
3	2/36
4	3/36
5	4/36
6	5/36
7	6/36
8	5/36
9	4/36
10	3/36
11	2/36
12	1/36

This is called the **probability mass function** (pmf) table of the random variable  $X$ . If  $S$  is the sample space, then  $P(S) = P(\text{union of all events of the form } X = x) = 1$  (verify from table).

# Random variable: Notation

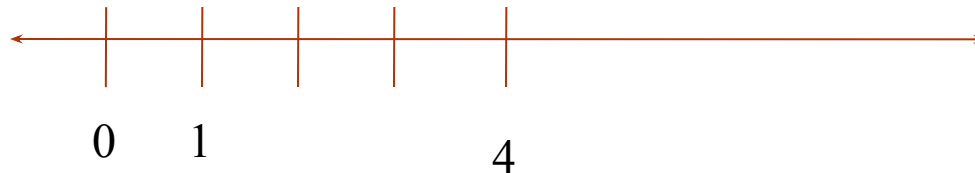
- A random variable is usually denoted by an *upper case* alphabet.
- Individual values the random variable can acquire are denoted by *lower case*.

# Random variable: discrete

- Random variables whose values can be written as a finite or infinite sequence are called **discrete random variables**.
- Example: results of coin toss or random dice experiments
- The probability that a random variable  $X$  takes on value  $x$ , i.e.  $P(X=x)$ , is called as the **probability mass function**.

# Random variable: continuous

- Random variables that can take on values within a continuum are called **continuous random variables**.
- Example: the dimensions (length, height, width, weight) of an object are usually continuous quantities, direction of a vector, amount of water that can be stored in a 4 litre jar is a continuous random variable in the interval  $[0,4]$ .



# Random variable: continuous

- For a continuous random variable, the probability that it takes on any *particular* value within a continuum is **zero!**
- Why? Because there are infinitely many values – say in the interval  $[0,4]$  in the example on the previous slide. Each value will be equally likely.
- **Note:** Zero probability in case of continuous random variables does *not* mean the event will *never* occur! This differs from the discrete case.



# Random variable: continuous

- Hence for a continuous random variable  $X$ , we consider the **cumulative distribution function** (cdf)  $F_X(x)$  defined as  $P\{X \leq x\}$ .
- The cdf is basically the probability that  $X$  takes on a value less than or equal to  $x$ .
- The cdf can be used to compute **cumulative interval measures**, that is the probability that  $X$  takes on a value greater than  $a$  and less than or equal to  $b$ , i.e.  $P(a < X \leq b) = F_X(b) - F_X(a)$ .

# Random variable: continuous - example

- Consider a cdf of the form:

$$F_X(x) = 0 \text{ for } x \leq 0, \text{ and}$$

$$F_X(x) = 1 - \exp(-x^2) \text{ otherwise}$$

- To find: probability that  $X$  exceeds 1

- $P(X > 1) = 1 - P(X \leq 1) = 1 - F_X(1) = e^{-1}$

# Probability Density Function (pdf)

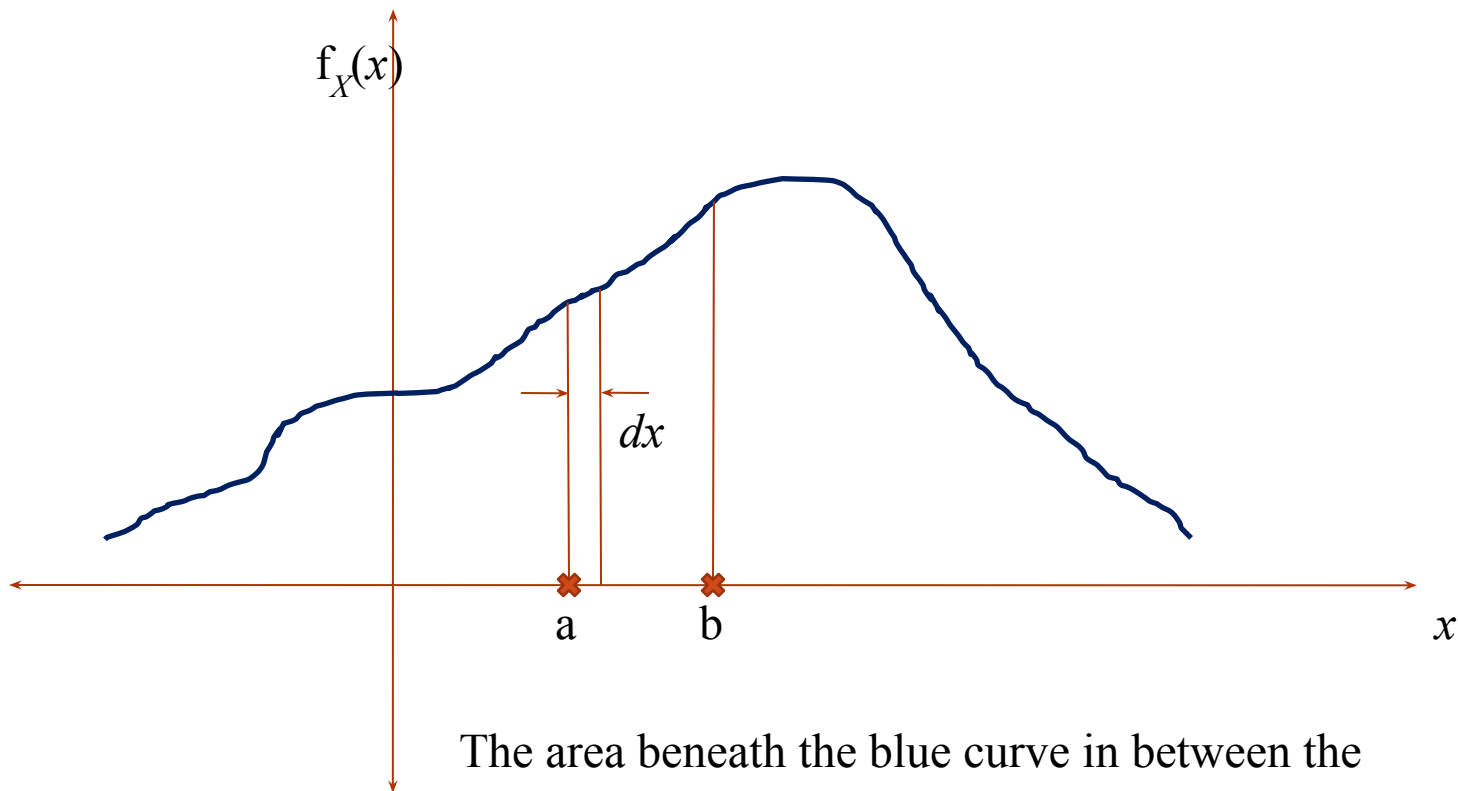
- The pdf of a random variable  $X$  at a value  $x$  is the *derivative* of its cumulative distribution function (cdf) at that value  $x$ .
- It is a non-negative function  $f_X(x)$  such that for any set  $B$  of real numbers, we have  $P\{X \in B\} = \int_B f_X(x)dx$

$$\int_{-\infty}^{+\infty} f_X(x)dx = 1$$

- Properties:

$$P(a \leq X \leq b) = \int_a^b f_X(x)dx = F_X(b) - F_X(a)$$

$$P(X = a) = \int_a^a f_X(x)dx = 0$$



The area beneath the blue curve in between the lines  $x = a$  and  $x = b$  is the **cumulative interval measure**  $P(a < X \leq b) = F_X(b) - F_X(a)$ .

$f_X(a)dx$  = probability that the random variable  $X$  takes on values between  $a$  and  $a+dx$ .

# Probability Density Function

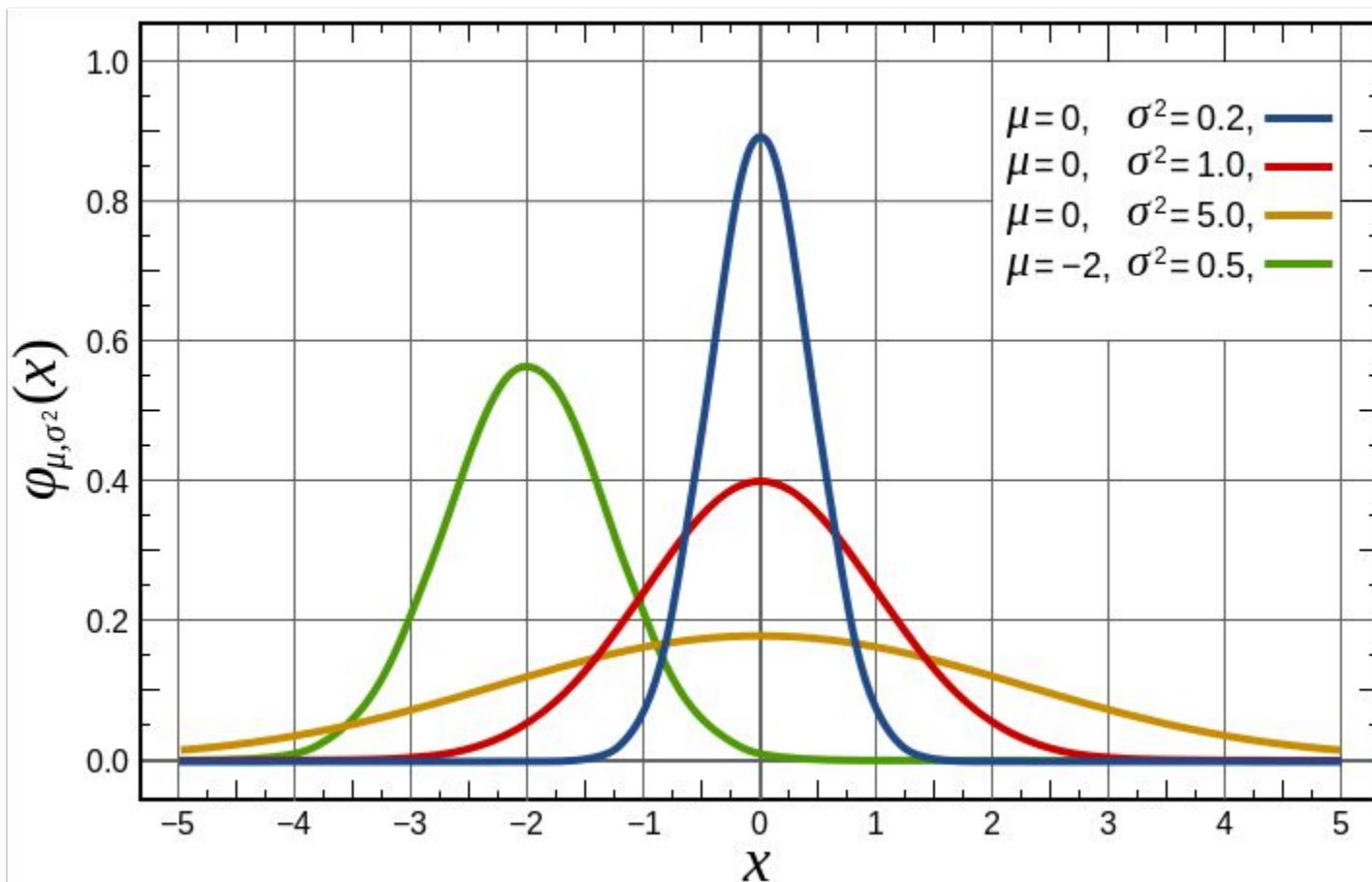
- Another way of looking at this concept:

$$P\{a - \varepsilon / 2 \leq X \leq a + \varepsilon / 2\} = \int_{a - \varepsilon / 2}^{a + \varepsilon / 2} f_X(x) dx \approx \varepsilon f(a)$$

$$f_X(a) = \lim_{\varepsilon \rightarrow 0} \frac{P\{a - \varepsilon / 2 \leq X \leq a + \varepsilon / 2\}}{\varepsilon}$$

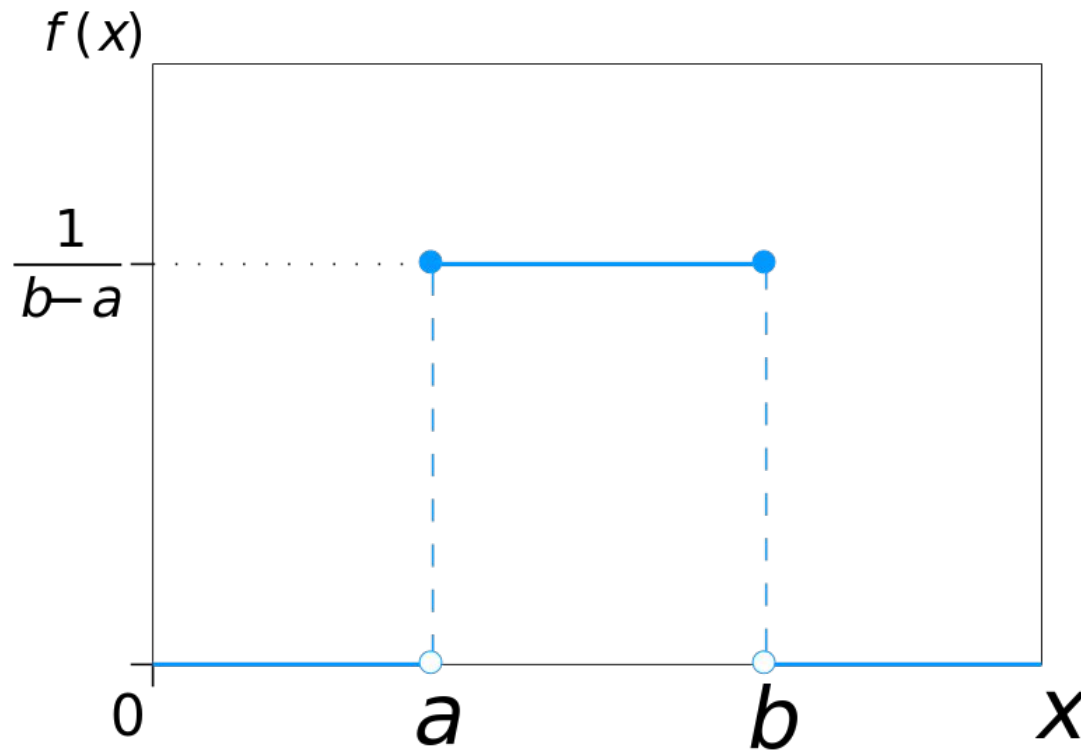
# Examples: Popular families of PDFs

- Gaussian (normal) pdf:  $f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/(2\sigma^2)}$



# Examples: Popular families of PDFs

- Bounded uniform pdf:  $f_X(x) = \frac{1}{(b-a)}, a \leq x \leq b$   
= 0 otherwise



# Expected Value (Expectation) of a random variable

- It is also called the **mean value** of the random variable.
- For a discrete random variable  $X$ , it is defined as:

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

- For a continuous random variable  $X$ , it is defined as:

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

- The expected value should not be (mis)interpreted to be the value that  $X$  *usually* takes on – it's the average value, *not* the “most frequently occurring value”.



# Expected Value (Expectation) of a random variable

- For some pdfs, the expected value is not always defined, i.e. the integral below may not have a finite value.

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

- One example is the pdf for the Pareto distribution (under some parameters) given as:

$$f_X(x | \alpha, x_m) = \frac{\alpha x_m^\alpha}{x^{\alpha+1}} \text{ for } x \geq x_m, \text{ otherwise } 0$$

$$E(X) = \alpha x_m^\alpha \left( \frac{x^{1-\alpha}}{1-\alpha} \right)_{x_m}^{\infty} = \infty \text{ if } \alpha < 1$$

$x_m$  and  $\alpha$  are parameters of the pdf for the Pareto distribution. Verify this result for  $E(X)$  on your own.

# Expected Value (Expectation) of a random variable

- Likewise for some discrete random variables which take on infinitely many values, the expected value may not be defined, i.e. we may have

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

- Example:

$$P(X = x) = k / x^2 \text{ for } x \geq 1, x \in \mathbb{Z}^+$$

$$E(X) = \sum_{x=1}^{\infty} x P(X = x) = k \sum_{x=1}^{\infty} 1/x = \infty \longrightarrow \text{See [here](#).}$$

$$\text{Note: } \sum_{x=1}^{\infty} P(X = x) = \sum_{x=1}^{\infty} k / x^2 = 1 \text{ if } k = 6 / \pi^2 \longrightarrow \text{See [here](#).}$$

# Expected Value: examples

- The expected value that shows up when you throw a die is  $1/6(1+2+3+4+5+6) = 3.5$ .
- The game of roulette consists of a ball and wheel with 38 numbered pockets on its side. The ball rolls and settles on one of the pockets. If the number in the pocket is the same as the one you guessed, you win \$35 (probability  $1/38$ ), otherwise you lose \$1 (probability  $37/38$ ). The expected value of the amount you earn after one trial is:  $(-1)37/38 + (35)1/38 = \$-0.0526$

## A Game of Roulette



[https://en.wikipedia.org/wiki/Roulette#/media/File:Roulette\\_casino.JPG](https://en.wikipedia.org/wiki/Roulette#/media/File:Roulette_casino.JPG)

# Expected value of a function of random variable

- Consider a function  $g(X)$  of a discrete random variable  $X$ . The expected value of  $g(X)$  is defined as (provided the summation is well-defined):

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

- For a continuous random variable, the expected value of  $g(X)$  is defined as (provided the integral is well-defined):

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

- This is called the Law of the Unconscious Statistician (LOTUS). It is something which requires proof, but is stated as if it were obvious.

# Proof: Law of the Unconscious Statistician

Let  $y = g(x)$ ,  $g$  is monotonic and differentiable.

$$g(g^{-1}(x)) = x, \therefore g'(g^{-1}(x))[g^{-1}]'(x) = 1$$

$$\therefore g'(g^{-1}(x)) = 1/[g^{-1}]'(x)$$

$$\text{Now } E(g(X)) = \int_{-\infty}^{\infty} g(x)f_X(x)dx$$

$$= \int_{-\infty}^{\infty} yf_X(g^{-1}(y))dx$$

$$= \int_{-\infty}^{\infty} yf_X(g^{-1}(y))dy / g'(g^{-1}(y))$$


# Proof (continued): Law of the Unconscious Statistician

$$E(g(X)) = \int_{-\infty}^{\infty} y f_X(g^{-1}(y)) dy / g'(g^{-1}(y))$$

$$\text{Now, } F_Y(y) = P(Y \leq y) = P(X \leq g^{-1}(y)) = F_X(g^{-1}(y))$$

Differentiating both sides, we have

$$f_Y(y) = f_X(g^{-1}(y)) / g'(g^{-1}(y))$$



This step assumes that  $g$  is a strictly increasing function. There is a fix to this, if  $g$  were strictly decreasing. What is it?

Substituting, we have

$$E(g(X)) = \int_{-\infty}^{\infty} y f_Y(y) dy$$

This justifies the previous formula for  $E(g(X))$ . Most textbooks just gloss over the formula, ignoring the fact that it requires proof.

# Properties of expected value

$$E(ag(X) + b) = \int_{-\infty}^{+\infty} (ag(x) + b)f_X(x)dx$$

$$= \int_{-\infty}^{+\infty} ag(x)f_X(x)dx + \int_{-\infty}^{+\infty} bf_X(x)dx$$

$$= aE(g(X)) + b \text{ --- why?}$$

This property is called the **linearity** of the expected value. In general, a function  $f(x)$  is said to be linear in  $x$  if  $f(ax+b) = af(x)+f(b)$  where  $a$  and  $b$  are constants. In this case, the expected value is not a function but an operator (it takes a function as input). An operator  $E$  is said to be linear if  $E(af(x) + b) = a E(f(x)) + E(b)$ . This is equal to  $aE(f(x)) + b$  for the expectation operator.



# Properties of expected value

- Consider a set of random variables  $X_1, X_2, \dots, X_n$ ; a set of functions  $g_1, g_2, \dots, g_n$ . Then we have:

$$E\left(\sum_{i=1}^n a_i g_i(X_i) + b_i\right) = \sum_{i=1}^n (a_i E[g_i(X_i)] + b_i)$$

$a_i, b_i$  are  
scalars

- This also forms a notion of the linearity of the expectation operator.
- Note: for a general nonlinear function  $g$ , we have:

$$E(g(X)) \neq g(E(X))$$

# Properties of expected value

Suppose you want to predict the value of a random variable with a known mean. On an average, what value will yield the least squared error?

Let  $X$  be the random variable and  $c$  be its predicted value.

We want to find  $c$  such that  $E((X - c)^2)$  is minimized.

Let  $\mu$  be the mean of  $X$ .

Then

$$\begin{aligned} E((X-c)^2) &= E((X - \mu + \mu - c)^2) \\ &= E((X - \mu)^2 + (\mu - c)^2 + 2(X - \mu)(\mu - c)) \\ &= E((X - \mu)^2) + E((\mu - c)^2) + 2E((X - \mu)(\mu - c)) \\ &= E((X - \mu)^2) + (\mu - c)^2 + 0 \\ &\geq E((X - \mu)^2) \end{aligned}$$

The expected value is the value that yields the least mean squared prediction error!

# The median

- What minimizes the following quantity?

$$J(c) = \int_{-\infty}^{+\infty} |x - c| f_X(x) dx$$

$$J(c) = \int_{-\infty}^c |x - c| f_X(x) dx + \int_c^{\infty} |x - c| f_X(x) dx$$

$$= \int_{-\infty}^c (c - x) f_X(x) dx + \int_c^{\infty} (x - c) f_X(x) dx$$

$$= \int_{-\infty}^c c f_X(x) dx - \int_{-\infty}^c x f_X(x) dx + \int_c^{\infty} x f_X(x) dx - \int_c^{\infty} c f_X(x) dx$$

$$= cF_X(c) - \int_{-\infty}^c x f_X(x) dx + \int_c^{\infty} x f_X(x) dx - c(1 - F_X(c))$$

# The median

$$J(c) = cF_X(c) - \int_{-\infty}^c xf_X(x)dx + \int_c^{\infty} xf_X(x)dx - c(1 - F_X(c))$$

$$J(c) = cF_X(c) - \int_{-\infty}^c q(x)dx + \int_c^{\infty} q(x)dx - c(1 - F_X(c))$$

$$q(x) = xf_X(x)$$

$$J(c) = cF_X(c) - (Q(c) - Q(-\infty)) + (Q(\infty) - Q(c)) - c(1 - F_X(c))$$

$$Q(c) = \int_{-\infty}^c xf_X(x)dx$$

$$= 2cF_X(c) - c - 2Q(c) + Q(\infty) + Q(-\infty)$$

In this derivation, we are assuming that the two definite integrals of  $q(x)$  exist! This proof won't go through otherwise.

# The median

$$J(c) = 2cF_X(c) - c - 2Q(c) + Q(\infty) + Q(-\infty)$$

$$J'(c) = 0$$

$$\therefore \underline{2cf_X(c) + 2F_X(c) - 1 - 2q(c)} = 0$$

$$\therefore 2cf_X(c) + 2F_X(c) - 1 - 2cf_X(c) = 0$$

$$\therefore 2F_X(c) - 1 = 0$$

$$\therefore F_X(c) = 1/2$$

This is the median – by definition and it minimizes  $J(c)$ . We can double check that  $J'(c) \geq 0$ . Notice the peculiar definition of the median for the continuous case here! This definition is not conceptually different from the discrete case, though. Also, note that the median will not be unique if  $F_X$  is not differentiable at  $c$ . This happens when  $F_X$  is not strictly increasing in some interval – say  $K = [c, c+\varepsilon]$  or  $[c-\varepsilon, c]$ . In such cases, all  $y \in K$  will qualify as medians and all of them will produce the same value of  $J(y)$ . This is because  $f_X(y) = 0$  for  $y \in K$ .

# Variance

- The **variance** of a random variable  $X$  tells you how much its values deviate from the mean – on an average.
- The definition of variance for a continuous r.v. with mean  $\mu$  is:

$$Var(X) = E[(X - \mu)^2] = \int_{-\infty}^{\infty} (x - \mu)^2 f_X(x) dx$$

- For a discrete r.v., the integration is replaced by a summation:

$$Var(X) = E[(X - \mu)^2] = \sum_i (x_i - \mu)^2 P(X = x_i)$$

- The positive square-root of the variance is called the **standard deviation**.
- Low-variance probability mass functions or probability densities tend to be concentrated around one point. High variance densities are spread out.

# Existence?

- For some distributions, the variance (and hence standard deviation) may not be defined, because the integral may not have a finite value.
- Example: Pareto distribution (see slides on expectation for definition) for  $\alpha < 2$ .
- Note in some cases the mean is defined, but the variance is not. In some cases both are undefined. However, if the mean is undefined, then the variance will be undefined too (why?).

# Variance: Alternative expression

- The definition of variance is:

$$Var(X) = E[(X - \mu)^2] = \int_{-\infty}^{\infty} (x - \mu)^2 f_X(x) dx$$

- Alternative expression:

$$\begin{aligned} Var(X) &= E[(X - \mu)^2] = E[X^2 + \mu^2 - 2X\mu] \\ &= E[X^2] + \mu^2 - 2E[X]\mu \\ &= E[X^2] + \mu^2 - 2\mu^2 \text{ --- why?} \\ &= E[X^2] - \mu^2 \\ &= E[X^2] - (E[X])^2 \end{aligned}$$



# Variance: properties

- Property:

$$\text{Var}(aX + b) = E[(aX + b - E(aX + b))^2]$$

$$= E[(aX + b - (a\mu + b))^2]$$

$$= E[a^2(X - \mu)^2]$$

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

# Probabilistic inequalities

- Sometimes we know the mean or variance of a random variable, and want to guess the probability that the random variable can take on a certain value.
- The exact probability can usually not be computed as the information is too less. But we can get upper or lower bounds on this probability which can influence our decision-making processes.

# Probabilistic inequalities

- Example: Let's say the average annual salary offered to a CSE Btech-4 student at IITB is \$100,000. What's the probability that you (i.e. a randomly chosen student) will get an offer of \$110,000 or more? **Additionally, if you were told that the variance of the salary was 50,000, what's the probability that your package is between \$90,000 and \$110,000?**

# Markov's inequality

- Let  $X$  be a random variable that takes only *non-negative* values. For any  $a > 0$ , we have

$$P\{X \geq a\} \leq E[X] / a$$

- Proof: next slide

# Markov's inequality

● Proof:

$$E[X] = \int_0^{\infty} x f_X(x) dx$$

$$= \int_0^a x f_X(x) dx + \int_a^{\infty} x f_X(x) dx$$

$$\geq \int_a^{\infty} x f_X(x) dx$$

$$\geq \int_a^{\infty} a f_X(x) dx$$

$$= a \int_a^{\infty} f_X(x) dx$$

$$= a P\{X \geq a\} \quad \therefore P\{X \geq a\} \leq E[X] / a$$

# Chebyshev's inequality

- For a random variable  $X$  with mean  $\mu$  and variance  $\sigma^2$ , we have for any value  $k > 0$ ,

$$P\{|X - \mu| \geq k\} \leq \frac{\sigma^2}{k^2}$$

- Proof: follows from Markov's inequality

$(X - \mu)^2$  is a non - negative random variable

$$\therefore P\{(X - \mu)^2 \geq k^2\} \leq E[(X - \mu)^2] / k^2 = \sigma^2 / k^2$$

$$\therefore P\{|X - \mu| \geq k\} \leq \sigma^2 / k^2$$

# Chebyshev's inequality: another form

- For a random variable  $X$  with mean  $\mu$  and variance  $\sigma^2$ , we have for any value  $k > 0$ ,

$$P\{|X - \mu| \geq k\} \leq \frac{\sigma^2}{k^2}$$

- If I replace  $k$  by  $k\sigma$ , I get the following:

$$P\{|X - \mu| \geq k\sigma\} \leq \frac{1}{k^2}$$

# Back to counting money! 😊

- Let  $X$  be the random variable indicating the annual salary offered to you when you reach Btech-4 😊
- Then

$$P\{X \geq 110K\} \leq \frac{100K}{110K} = 0.9090 \approx 90\%$$

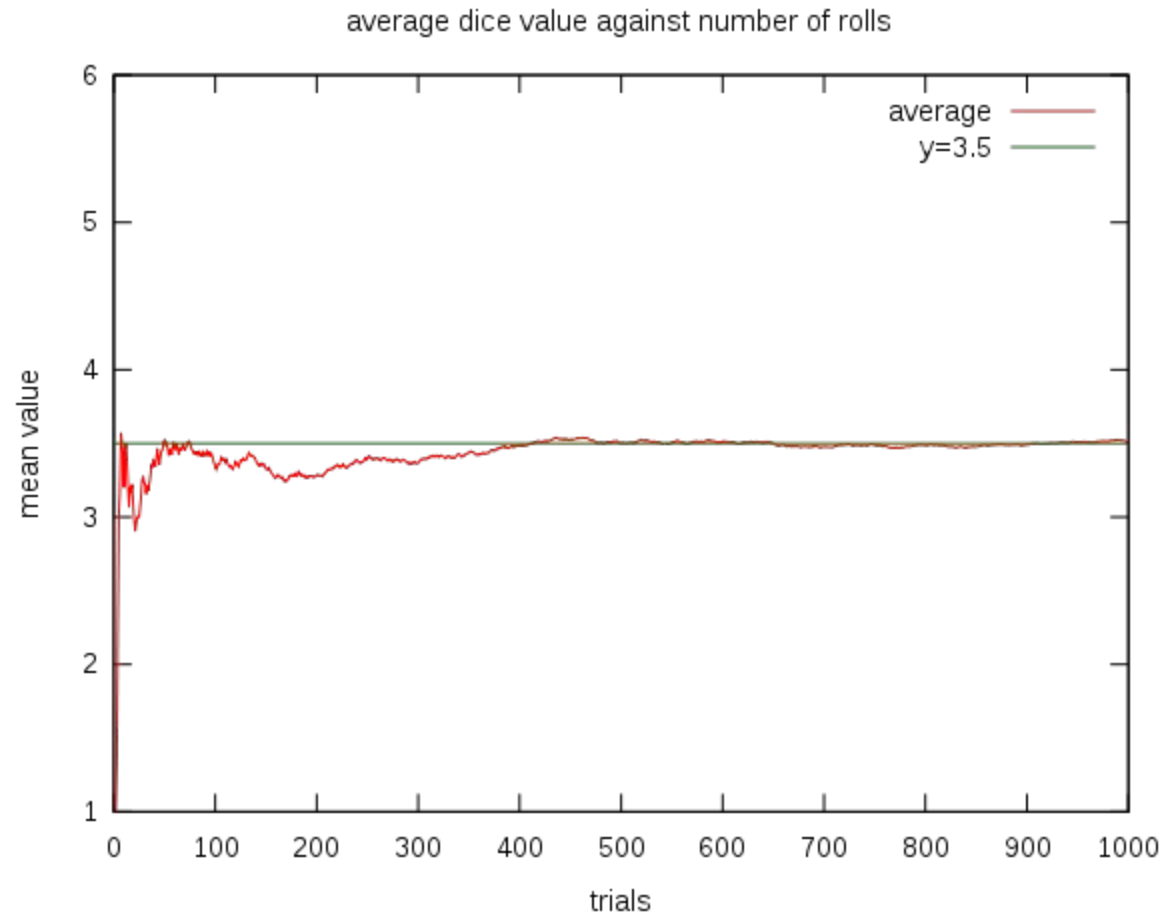
$$P\{|X - 100K| \geq 10K\} \leq \frac{50K}{10K \times 10K} = 0.0005 \approx 0.05\%$$

$$\therefore P\{|X - 100K| < 10K\} = 1 - 0.05\% = 99.5\%$$



# Back to the expected value

- When I tell you that the expected value of a random die variable is 3.5, what does this mean?
- If I throw the die  $n$  times, and average the results, I should get a value close to 3.5 provided  $n$  is very large (not valid if  $n$  is small).
- As  $n$  increases, the average value should move closer and closer towards 3.5.
- That's our basic intuition!



[https://en.wikipedia.org/wiki/Law\\_of\\_large\\_numbers](https://en.wikipedia.org/wiki/Law_of_large_numbers)

# Back to the expected value: weak law of large numbers

- This intuition has a rigorous theoretical justification in a theorem known as the **weak law of large numbers**.
- Let  $X_1, X_2, \dots, X_n$  be a sequence of independent and identically distributed random variables each having mean  $\mu$ . Then for any  $\varepsilon > 0$ , we have:

$$P\left\{\left|\frac{X_1 + X_2 + \dots + X_n}{n} - \mu\right| > \varepsilon\right\} \rightarrow 0 \text{ as } n \rightarrow \infty$$

# Back to the expected value: weak law of large numbers

- Let  $X_1, X_2, \dots, X_n$  be a sequence of independent and identically distributed random variables each having mean  $\mu$ . Then for any  $\varepsilon > 0$ , we have:

$$P\left\{\left|\frac{X_1 + X_2 + \dots + X_n}{n} - \mu\right| > \varepsilon\right\} \rightarrow 0 \text{ as } n \rightarrow \infty$$

Empirical (or sample) mean

- Proof: follows immediately from Chebyshev's inequality

$$E\left(\frac{X_1 + X_2 + \dots + X_n}{n}\right) = \mu, \text{Var}\left(\frac{X_1 + X_2 + \dots + X_n}{n}\right) = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n},$$

$$\therefore P\left\{\left|\frac{X_1 + X_2 + \dots + X_n}{n} - \mu\right| > \varepsilon\right\} \leq \frac{\sigma^2}{n\varepsilon^2}$$

$$\therefore \lim_{n \rightarrow \infty} P\left\{\left|\frac{X_1 + X_2 + \dots + X_n}{n} - \mu\right| > \varepsilon\right\} = 0$$

# Comments on weak law of large numbers

- The assumption of “identically distributed” is not strictly necessary (even though it is stated that way in many texts), but the random variables should have the same mean  $\mu$ .
- The previous proof assumes they have the same variance – if not the  $n\sigma^2$  term in the RHS numerator would be replaced by  $n$  times the average variance of the random variables. The proof as such still goes through.
- The independence assumption is not strictly necessary – the random variables just need to be *pair-wise uncorrelated* (more on this in later slides) – so that the variances can add up as shown in the proof.
- The law assumes that the random variables have a well-defined expected value and variance. Otherwise, the law may not hold.

# The strong law of large numbers

- The strong law of large numbers states the following:

$$P(\lim_{n \rightarrow \infty} \frac{X_1 + X_2 + \dots + X_n}{n} = \mu) = 1$$

- This is stronger than the weak law because this states that the probability of the desired event (that the empirical mean is equal to the actual mean) is **equal to 1** given enough samples. The weak laws states that it **tends to 1** given enough samples.
- The proof of the strong law is formidable and beyond the scope of our course.

# (The incorrect) Law of averages

- As laymen we tend to believe that if something has been going wrong for quite some time, it will suddenly turn right – using the law of averages.
- This supposed law is actually a fallacy – it reflects wishful thinking, and the core mistake is that we mistake the distribution of samples among a small set of outcomes for the distribution of a larger set.
- This is also called as **Gambler's fallacy**.

# (The incorrect) Law of averages

- Let's say a gambler *independently* tosses an unbiased coin 20 times, and gets a head each time. He now applies the “law of averages” and believes that it is more likely that the next coin toss will yield a tail.
- The mistake is as follows: The probability of getting all 21 heads =  $(1/2)^{21}$ . The probability of getting 20 heads and 1 tail also =  $(1/2)^{21}$ .



# Joint distributions/pdfs/pmfs

# Jointly distributed random variables

- Many times in statistics, one needs to model relationships between two or more random variables – for example, your CPI at IITB and the annual salary offered to you during placements!
- Another example: average amount of sugar consumed per day and blood sugar level recorded in a blood test.
- Another example: literacy level and crime rate.

# Joint CDFs

- Given continuous random variables  $X$  and  $Y$ , their **joint cumulative distribution function** (cdf) is defined as:

$$F_{XY}(x, y) = P(X \leq x, Y \leq y)$$

- The distribution of either random variable (called as marginal cdf) can be obtained from the joint distribution as follows:

$$F_X(x) = P(X \leq x, Y \leq \infty) = F_{XY}(x, \infty)$$

$$F_Y(y) = P(X \leq \infty, Y \leq y) = F_{XY}(\infty, y)$$

I'll explain this  
a few slides  
further down

- These definitions can be extended to handle more than two random variables as well.

# Joint PMFs

- Given two discrete random variables  $X$  and  $Y$ , their **joint probability mass function** (pmf) is defined as:

$$p_{XY}(x_i, y_j) = P(X = x_i, Y = y_j)$$

- The pmf of either random variable (called as marginal pmf) can be obtained from the joint distribution as follows:

$$P\{X = x_i\} = P\left(\bigcup_j \{X = x_i, Y = y_j\}\right)$$

$$= \sum_j P\{X = x_i, Y = y_j\} = \sum_j p(x_i, y_j)$$

Why?

# Joint PMFs: Example

- Consider that in a city 15% of the families are childless, 20% have only one child, 35% have two children and 30% have three children. Let us suppose that male and female child are equally likely and independent.
- What is the probability that a randomly chosen family has no children?
- $P(B = 0, G = 0) = 0.15 = P(\text{no children})$
- Has 1 girl child and no boy child?
- $P(B=0, G=1) = P(1 \text{ child}) P(G=1 | 1 \text{ child}) = 0.2 \times 0.5 = 0.1$
- Has 3 girls?
- $P(B = 0, G = 3) = P(3 \text{ children}) P(G=3 | 3 \text{ Children}) = 0.3 \times (0.5)^3$
- Has 2 boys and 1 girl?
- $P(B = 2, G = 1) = P(3 \text{ children}) P(B = 2, G = 1 | 3 \text{ children}) = 0.3 \times (1/8) \times 3 = 0.1125$  (all 8 combinations of 3 children are equally likely. Out of these there are 3 of the form 2 boys + 1 girl)

# Joint PDFs

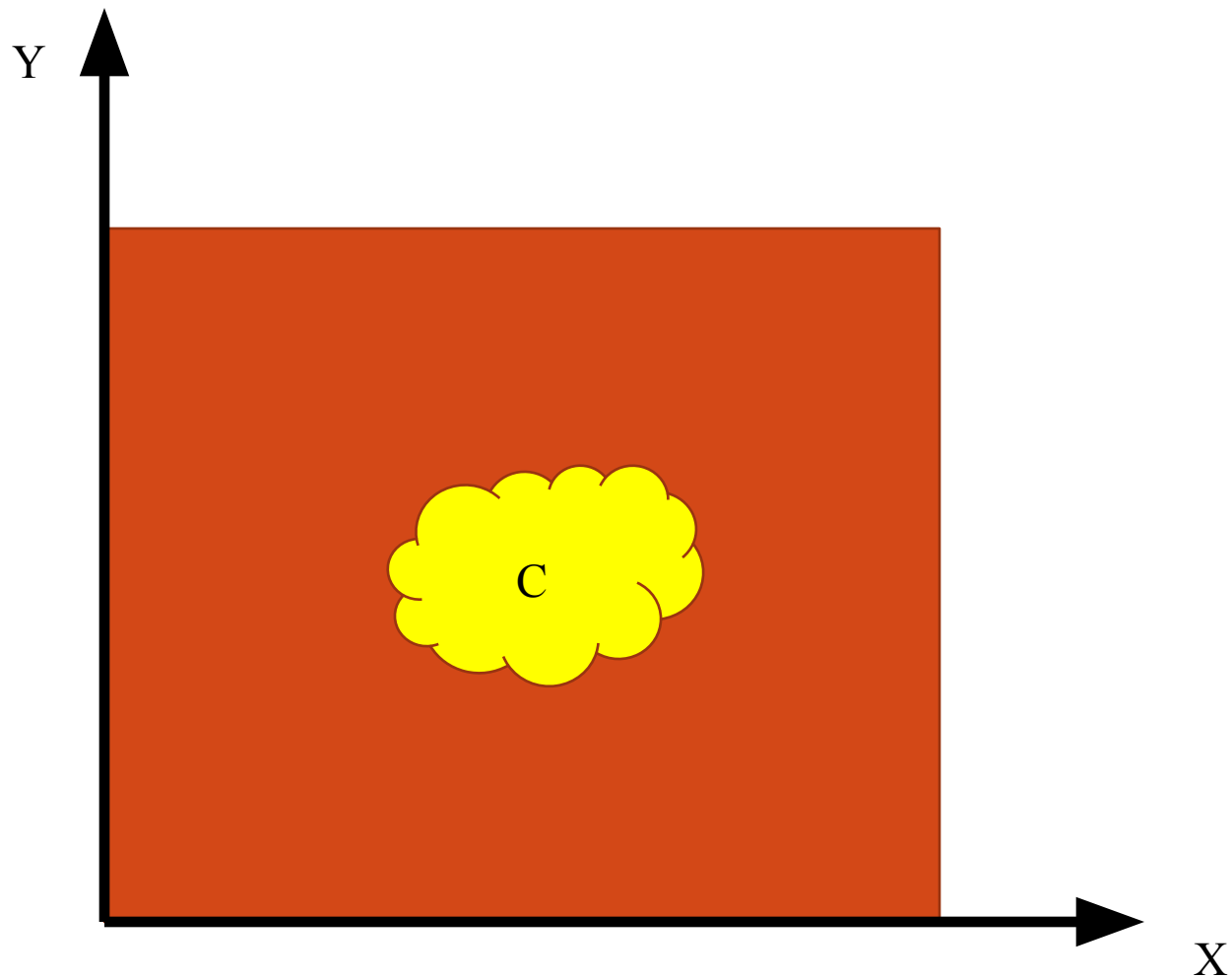
- For two jointly continuous random variables  $X$  and  $Y$ , the joint pdf is a non-negative function  $f_{XY}(x,y)$  such that for *any* set  $C$  in the two-dimensional plane, we have:

$$P\{(X,Y) \in C\} = \iint_{(x,y) \in C} f_{XY}(x,y) dx dy$$

- The joint CDF can be obtained from the joint PDF as follows:

$$F_{XY}(a,b) = \int_{-\infty}^a \int_{-\infty}^b f_{XY}(x,y) dx dy$$

$$f_{XY}(a,b) = \frac{\partial^2}{\partial x \partial y} F_{XY}(x,y) \big|_{x=a,y=b}$$



The joint probability that  $(X,Y)$  belongs to any arbitrary-shaped region in the  $XY$ -plane is obtained by integrating the joint pdf of  $(X,Y)$  over that region (eg: region  $C$ )

# Joint and marginal PDFs

- The **marginal pdf** of a random variable can be obtained by integrating the joint pdf w.r.t. the other random variable(s):

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

$$f_Y(y) = \int_{-\infty}^{\infty} f_{XY}(x, y) dx$$

$$f_Y(y) = \int_{-\infty}^{\infty} f_{XY}(x, y) dx$$

$$\int_{-\infty}^a f_X(x) dx = \int_{-\infty}^a \int_{-\infty}^{\infty} f_{XY}(x, y) dy dx$$

$$\therefore F_X(a) = F_{XY}(a, \infty)$$



# Independent random variables

- Two continuous random variables are said to be **independent** if and only if:

$$\forall x, \forall y, f_{XY}(x, y) = f_X(x)f_Y(y)$$

i.e., the joint pdf is equal to the product of the marginal pdfs.

- For independent random variables, the joint CDF is also equal to the product of the marginal CDFs:

$$F_{XY}(x, y) = F_X(x)F_Y(y) \quad \text{Try proving this yourself!}$$

# Independent random variables

- Some  $n$  continuous random variables  $X_1, X_2, \dots, X_n$  are said to be **mutually independent** if and only if for any finite subset of  $k$  random variables  $X_{i1}, X_{i2}, \dots, X_{ik}$  and finite sequence of number  $x_1, x_2, \dots, x_k$ , the events  $X_{i1} \leq x_1, X_{i2} \leq x_2, \dots, X_{ik} \leq x_k$  are mutually independent.

- As a consequence

$$\forall x_1, x_2, \dots, x_n,$$

$$f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n) = f_{X_1}(x_1) f_{X_2}(x_2) \dots f_{X_n}(x_n)$$

i.e., the joint pdf is equal to the product of all  $n$  marginal pdfs.

- Note that this condition is *stronger* than pairwise independence!

$$\forall (x_i, x_j), 1 \leq i \leq n, 1 \leq j \leq n, i \neq j,$$

$$f_{X_i, X_j}(x_i, x_j) = f_{X_i}(x_i) f_{X_j}(x_j)$$

# Independent random variables

- Mutual independence between  $n$  random variables implies that they are pairwise independent, or in fact,  $k$ -wise independent for any  $k < n$ .
- But pairwise independence does not necessarily imply mutual independence.
- Example: Consider a sample space  $\{1,2,3,4\}$  where each singleton element is equally likely to be chosen.

# Independent random variables

- Consider  $A = \{1,2\}$ ,  $B = \{1,3\}$ ,  $C = \{1,4\}$ .
- Then  $P(A) = P(B) = P(C) = 1/2$ .  $P(ABC) = P(\{1\}) = 1/4 \neq P(A)P(B)P(C)$  implying that  $A, B, C$  are not mutually independent.
- But  $P(AB) = 1/4 = P(A)P(B)$  and likewise for  $AC$ ,  $BC$ .
- Let us define the random variables  $E_A$ ,  $E_B$ ,  $E_C$  which acquire the value 1 if events  $A, B, C$  respectively occur, and 0 otherwise.
- We can see that  $E_A$ ,  $E_B$ ,  $E_C$  are pair-wise independent random variables, but not mutually independent.

# Concept of covariance

- The covariance of two random variables  $X$  and  $Y$  is defined as follows:

$$\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$$

- Further expansion:

$$\begin{aligned}\text{Cov}(X, Y) &= E[(X - \mu_X)(Y - \mu_Y)] \\ &= E[XY - \mu_X Y - \mu_Y X + \mu_X \mu_Y] \\ &= E[XY] - \mu_X \mu_Y - \mu_Y \mu_X + \mu_X \mu_Y \text{ --- why?} \\ &= E[XY] - \mu_X \mu_Y \\ &= E[XY] - E[X]E[Y]\end{aligned}$$

# Concept of covariance: properties

- $\text{Cov}(X, Y) = \text{Cov}(Y, X)$
- $\text{Cov}(X, X) = \text{Var}(X)$  [verify this yourself!]
- $\text{Cov}(aX, Y) = a\text{Cov}(X, Y)$  [prove this!]
- Relationship with correlation coefficient:

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

# Concept of covariance: properties

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

*Proof :*

$$\begin{aligned}\text{Cov}(X + Z, Y) &= E[(X + Z)Y] - E[X + Z]E[Y] \\ &= E[XY + ZY] - E[X]E[Y] - E[Z]E[Y] \\ &= E[XY] - E[X]E[Y] + E[ZY] - E[Z]E[Y] \\ &= \text{Cov}(X, Y) + \text{Cov}(Z, Y)\end{aligned}$$

$$\text{Cov}\left(\sum_i X_i, Y\right) = \sum_i \text{Cov}(X_i, Y)$$

$$\text{Cov}\left(\sum_i X_i, \sum_j Y_j\right) = \sum_i \sum_j \text{Cov}(X_i, Y_j)$$

Try proving this yourself! Along similar lines as the previous one.

# Concept of covariance: properties

$$\text{Cov}(\sum_i X_i, Y) = \sum_i \text{Cov}(X_i, Y)$$

$$\text{Cov}(\sum_i X_i, \sum_j Y_j) = \sum_i \sum_j \text{Cov}(X_i, Y_j)$$

$$\begin{aligned}\text{Var}(\sum_i X_i) &= \text{Cov}(\sum_i X_i, \sum_i X_i) \\ &= \sum_i \sum_j \text{Cov}(X_i, X_j) \\ &= \sum_i \text{Cov}(X_i, X_i) + \sum_i \sum_{j \neq i} \text{Cov}(X_i, X_j) \\ &= \sum_i \text{Var}(X_i) + \sum_i \sum_{j \neq i} \text{Cov}(X_i, X_j)\end{aligned}$$

Notice that the variance of the sum of random variables is not equal to the sum of their individual variances. This is quite unlike the mean!



# Concept of covariance: properties

- For independent random variables  $X$  and  $Y$ ,  $\text{Cov}(X, Y) = 0$ , i.e.  $E[XY] = E[X]E[Y]$ .

- Proof:

$$\begin{aligned} E[XY] &= \sum_i \sum_j x_i y_j P\{X = x_i, Y = y_j\} \\ &= \sum_i \sum_j x_i y_j P\{X = x_i\} P\{Y = y_j\} \\ &= \sum_i x_i P\{X = x_i\} \sum_j y_j P\{Y = y_j\} \\ &= E[X]E[Y] \end{aligned}$$

$$\begin{aligned} \text{Cov}(X, Y) &= E[(X - \mu_X)(Y - \mu_Y)] \\ &= E[XY] - \mu_X E[Y] - \mu_Y E[X] + \mu_X \mu_Y \\ &= E[XY] - E[X]E[Y] = 0 \end{aligned}$$

# Concept of covariance: properties

- Given random variables  $X$  and  $Y$ ,  $\text{Cov}(X, Y) = 0$  does not necessarily imply that  $X$  and  $Y$  are independent!
- Proof: Construct a counter-example yourself!



# Conditional pdf/cdf/pmf

- Given random variables  $X$  and  $Y$  with joint pdf  $f_{XY}(x,y)$ , then the conditional pdf of  $X$  given  $Y = y$  is defined as follows:

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

- Conditional cdf  $F_{X|Y}(x,y)$ :

$$\begin{aligned} F_{X|Y}(x | y) &= \lim_{\delta \rightarrow 0} P(X \leq x | y \leq Y \leq y + \delta) = \int_{-\infty}^x f_{X|Y}(z | y) dz \\ &= \int_{-\infty}^x \frac{f_{X,Y}(z, y)}{f_Y(y)} dz \end{aligned}$$

# Conditional pdf/cdf/pmf

- Conditional cdf  $F_{X|Y}(x,y)$ :

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

$$= \frac{F_{XY}(x, y + \delta) - F_{XY}(x, y)}{(F_Y(y + \delta) - F_Y(y))}$$

$$\approx \frac{(\partial F_{XY}(x, y) / \partial y) \delta}{f_Y(y) \delta}$$

$$= \frac{(\partial F_{XY}(x, y) / \partial y)}{f_Y(y)}$$

$$\begin{aligned} f_{X|Y}(x | y) &= \frac{\partial}{\partial x} (F_{X|Y}(x | y)) = \frac{\partial}{\partial x} \left( \frac{\partial F_{XY}(x, y) / \partial y}{f_Y(y)} \right) \int_{-\infty}^{\infty} f_{X|Y}(x | y) dx = \int_{-\infty}^{\infty} \frac{f_{XY}(x, y)}{f_Y(y)} dx \\ &= \frac{\partial^2 F_{XY}(x, y) / \partial x \partial y}{f_Y(y)} = \frac{f_{XY}(x, y)}{f_Y(y)} = \frac{1}{f_Y(y)} \int_{-\infty}^{\infty} f_{XY}(x, y) dx = \frac{f_Y(y)}{f_Y(y)} = 1 \end{aligned}$$

# Conditional mean and variance

- Conditional densities or distributions can be used to define the **conditional mean** (also called **conditional expectation**) or **conditional variance** as follows:

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

$$Var(X | Y = y) = \int_{-\infty}^{\infty} (x - E(X | Y = y))^2 f_{X|Y}(x | y) dx$$

# Example

$$f(x, y) = 2.4x(2 - x - y), 0 < x < 1, 0 < y < 1$$

$= 0$  otherwise

Find conditional density of  $X$  given  $Y = y$ .

Find conditional mean of  $X$  given  $Y = y$ .

# Moment Generating Functions

# Definition

- The **moment** of random variable  $X$  of order  $n$  is defined as follows  $m_n = E(X^n)$ .
- The moment generating function (MGF) of a random variable  $X$  is defined as follows:

$$\phi_X(t) = E(e^{tX}) = \sum_x e^{tx} P(X = x) \text{ (discrete r.v.)}$$

$$= \int_{-\infty}^{\infty} f_X(x) e^{tx} dx \text{ (continuous r.v.)}$$




# Why is it so called?

- Because of:

$$e^{tX} = 1 + tX + (tX)^2 / 2! + (tX)^3 / 3! + \dots$$

$$\phi_X(t) = E(e^{tX}) = 1 + tm_1 + t^2 m_2 / 2! + t^3 m_3 / 3! + \dots$$


$$m_i = E(X^i), i \geq 1$$

# Key property

- Differentiating the MGF w.r.t. the parameter  $t$  yields the different moments of  $X$ .

$$\phi'_X(t) = \frac{d}{dt} \left( E(e^{tX}) \right) = E \left( \frac{d}{dt} (e^{tX}) \right) = E(Xe^{tX})$$

$$\phi'_X(0) = E(X)$$

$$\phi_X^{(2)}(t) = \frac{d}{dt} \left( E(Xe^{tX}) \right) = E(X^2 e^{tX})$$

$$\phi_X^{(2)}(0) = E(X^2)$$

....

$$\phi_X^{(n)}(0) = E(X^n)$$

# Other properties

- If  $Y = aX + b$ , then we have:  $\phi_Y(t) = e^{tb} \phi_X(at)$
- If  $Y$  and  $X$  are independent, then:  $\phi_{X+Y}(t) = \phi_X(t) \phi_Y(t)$
- Let  $X$  and  $Y$  be random variables. Let  $Z$  be a third r.v. which is equal to  $X$  with probability  $p$ , and equal to  $Y$  with probability  $1-p$ . Then we have:

$$\phi_Z(t) = p\phi_X(t) + (1-p)\phi_Y(t)$$

# Uniqueness

- For a discrete random variable with finite range, the MGF and PMF uniquely determine each other.
- Proof:

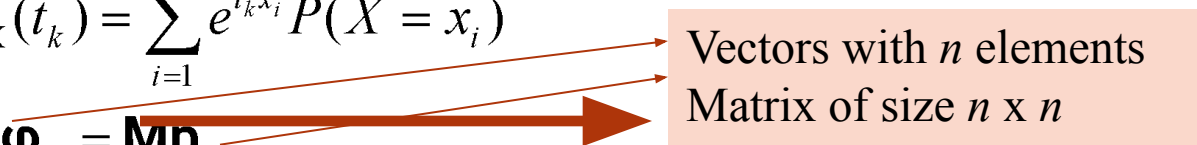
$$\phi_X(t) = E(e^{tX}) = \sum_x p(X=x)e^{tX} \therefore \text{PMF uniquely determines MGF.}$$

To prove the converse, consider that  $X$  takes on some  $n$  values.

Consider some  $n$  values of  $t$  as well. Then we have :

$$\phi_X(t_k) = \sum_{i=1}^n e^{t_k x_i} P(X = x_i)$$

$\therefore \boldsymbol{\phi}_X = \mathbf{M}\mathbf{p}$



Vectors with  $n$  elements  
Matrix of size  $n \times n$

The matrix  $\mathbf{M}$  has a special form that makes it invertible.

Hence  $\mathbf{p} = \mathbf{M}^{-1}\boldsymbol{\phi}_X$  is uniquely determined.

Proof here and here.

# Uniqueness: Another proof

- If two discrete random variables  $X$  and  $Y$  have MGFs  $\phi_X(t)$  and  $\phi_Y(t)$  that both exist and  $\phi_X(t) = \phi_Y(t)$  for all  $t$ , then  $X$  and  $Y$  have the same probability mass function.
- Proof for discrete random variables:

$$\phi_X(t) = \phi_Y(t)$$

$$\therefore \sum_x e^{tx} p(X=x) = \sum_y e^{ty} p(Y=y) = \sum_x e^{tx} p(Y=x)$$

$$\therefore \sum_x e^{tx} (p(X=x) - p(Y=x)) = 0$$

$$\therefore \sum_x s^x c_x = 0 \text{ where } s = e^t, c_x = p(X=x) - p(Y=x)$$

This is a polynomial in  $s$  with coefficients  $\{c_x\}$ . The polynomial can be 0 for all values of  $s$ , iff  $c_x = 0$ . Hence  $p(X=x) = p(Y=x)$  for all  $x$ .

# Uniqueness: Continuous case

- The uniqueness theorem is also applicable to continuous random variables, although we do not prove it here.