

Chapter 1

Background material

1.1 Distribution Summary

1. Binomial(n, p); $f(x) = \binom{n}{x} p^x (1-p)^{n-x}$ $x = 0, 1, \dots, n$

where $\binom{n}{x} = \frac{n!}{(x!)(n-x)!}$ *orange*

Consider n independent repetitions of an experiment each of which has only two possible outcomes, say (S, F) where

$P\{S\} = p$ is constant, i.e. the same for each experiment

Let $X = \#S$'s in n trials

Then $X \sim \text{Binomial}(n, p)$.

Example: Let X be the number of heads in n tosses of a fair coin.

Note: We often use the Binomial Distribution when we **Sample with Replacement**.

$$E(X) = np, \quad \text{Var}(X) = np(1-p)$$

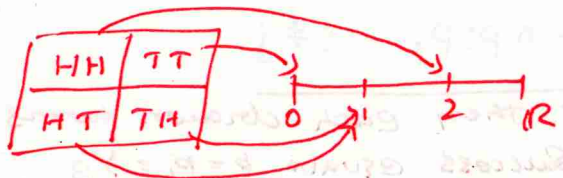
Since, x is the # heads in 2 tosses

Let, $n=2$. What is the Sample space of the experiment i.e., the domain of x ?

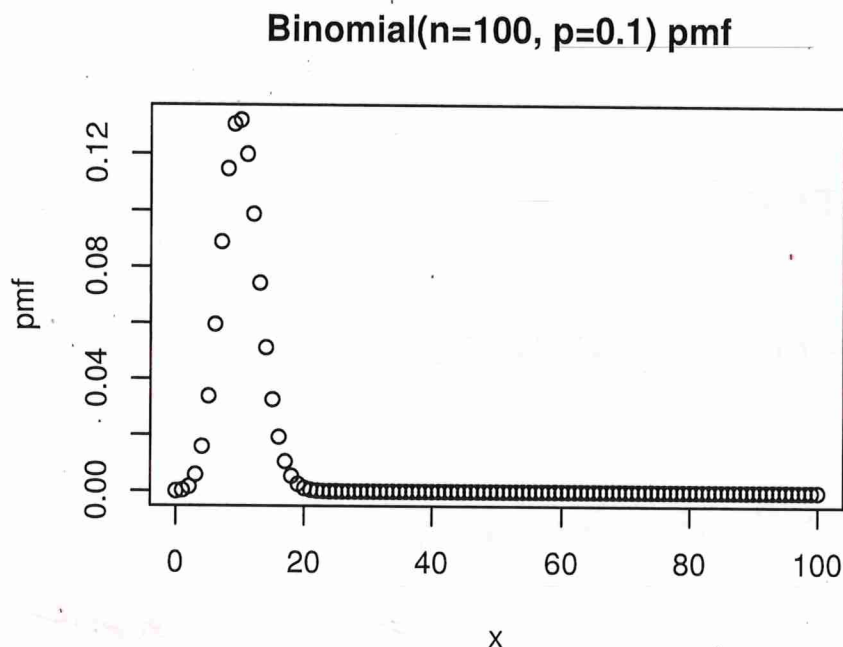
$$X = \{0, 1, 2\}$$

1

$$X \sim \text{Bin}(2, p=0.5)$$



$$\Omega = \{HH, HT, TH, TT\}$$



x is discrete

Figure 1.1: Binomial probability mass function, $n=100$, $p=0.1$

2. Multinomial(n, p_1, \dots, p_k); $f(x_1, \dots, x_k) = \binom{n}{x_1 \dots x_k} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}$

where $\binom{n}{x_1 \dots x_k} = \frac{n!}{(x_1!)(x_2!) \dots (x_k!)}$ and

$x_i = 0, 1, \dots, n$, such that $x_1 + \dots + x_k = n$ and $p_1 + \dots + p_k = 1$.



Consider n independent repetitions of an experiment for which each outcome can be classified in exactly one of k mutually exclusive ways, A_1, A_2, \dots, A_k .

Let

When $k=2$; this is Binomial(n, p_1)

$p_i = P\{\text{an outcome of one trial is of class } A_i\}$

$X_i = \# \text{ outcomes that are of class } i \text{ out of } n \text{ repetitions}$

X_i 's have marginal Binomial Distributions

Then $(X_1, X_2, \dots, X_k) \sim \text{Multinomial}(n, p_1, \dots, p_k)$

$$E(X_i) = np_i, \quad \text{Var}(X_i) = np_i(1-p_i)$$

$$\text{Cov}(X_i, X_j) = -np_i p_j \quad i \neq j$$

We draw marbles from a bag

O.R.

*2 B R
0 B 4*

*Let, $X_1 = \# \text{ red}$,
 $X_2 = \# \text{ Blue}$,
 $X_3 = \# \text{ Green Marbles}$*

$(X_1, X_2, X_3) \sim \text{Multi}(3, \frac{1}{3}, \frac{1}{3}, \frac{1}{3})$

Now, let's focus on only on red marbles; then each drawn turns into a yes/no trial, with a probability of success equals $p = p_1 = \frac{1}{3}$ and failure being $1-p = 1-p_1 = \frac{2}{3} = \frac{1}{3} + \frac{1}{3} = p_2 + p_3$.

$\therefore X_1 \sim \text{Bin}(3, \frac{1}{3})$ and similarly, $X_2 \sim \text{Bin}(2, \frac{1}{3})$, $X_3 \sim \text{Bin}(3, \frac{1}{3})$

Note: $\sum_{i=1}^k p_i = 1 \quad \sum_{i=1}^k X_i = n$

Example: Toss a fair die $n = 100$ times and let (X_1, X_2, \dots, X_6) be the observed frequencies of the numbers 1, 2, 3, 4, 5, 6 from the tosses of the die. Since the die is fair, then $p_i = 1/6$ for $i = 1, \dots, 6$.

3. **Negative Binomial** (r, p) ; $f(x) = \binom{x+r-1}{r-1} p^r (1-p)^x$, $x = 0, 1, \dots$

Consider independent repetitions of an experiment each of which, has exactly two possible outcomes, say (S, F) .

Let $P(S) = p$ constant, i.e. the same for each experiment

Let $X = \# F$'s before the r^{th} S

Then $X \sim \text{NegBin}(r, p)$

X is discrete

Example: Continue flipping a fair coin and stop when you observe the first head. $X =$ the number of tails before the first head has a Negative Binomial distribution with $r = 1$.

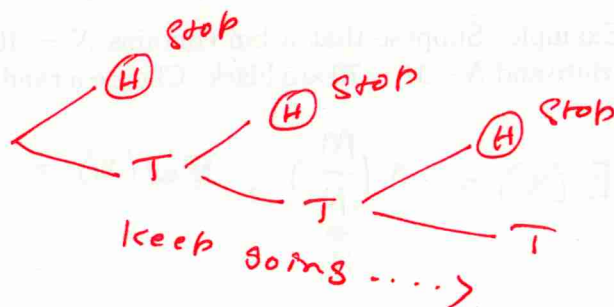
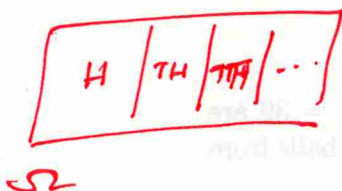
$$E(X) = \frac{r(1-p)}{p}, \quad \text{Var}(X) = \frac{r(1-p)}{p^2}$$

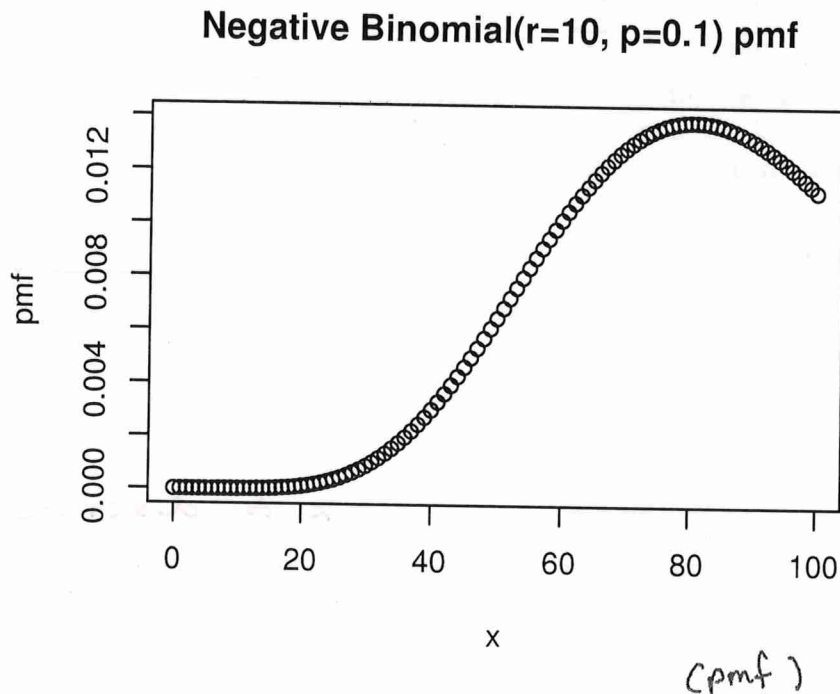
What is the sample space here?

$$\Omega = \{H, TH, TTH, TTTH, \dots\}$$

How many possible outcomes?

Countably infinitely many outcomes.



Figure 1.2: Negative Binomial probability mass function, $r=10, p=0.1$

4. **Geometric(p)**; is the same as Negative Binomial ($r=1, p$) $f(x) = p(1-p)^x; x=0, 1, 2, \dots$
 or, $f(x) = p(1-p)^{x-1}; x=1, 2, \dots$
5. **Hypergeometric(N, M, n)**; $f(x) = \frac{\binom{M}{x} \binom{N-M}{n-x}}{\binom{N}{n}}$ where $\max(0, n-M) \leq x \leq \min(n, M)$.

Consider a finite population of size N . Let each object in the population be characterized as either a S or F, where there are $M \leq N$ S's in the population. Draw a random sample of size n from the population without replacement.

Let $X = \# \text{ S's in the sample of size } n$.

Then $X \sim \text{Hypergeometric}(N, M, n)$.

Example: Suppose that a bin contains $N = 100$ balls, of which $M = 30$ are white and $N - M = 70$ are black. Choose a random sample of $n = 10$ balls from

$$E(X) = n \underbrace{\left(\frac{M}{N}\right)}_{\substack{\uparrow \\ p}}, \quad \text{Var}(X) = \frac{N-n}{N-1} n \underbrace{\left(\frac{M}{N}\right)}_{\substack{\uparrow \\ p}} \left[1 - \underbrace{\left(\frac{M}{N}\right)}_{\substack{\uparrow \\ p}}\right]$$

Explained at the end of this pdf
(last page).

the bin without replacement. X = the number of white balls in the sample has a Hypergeometric($N = 100, M = 30, n = 10$) distribution.



Example: A shipping container contains $N = 10,000$ iPhone 7's of which $M = 30$ are defective and the remainder are not defective. Choose a random sample of $n = 100$ iPhone 7's from the shipping container without replacement. Then X = the number of defectives in the sample has a Hypergeometric($N = 10,000, M = 30, n = 100$) distribution.

In this example, $n/N = 100/10,000 = 0.01 \leq 0.05$. Then X = the number of defectives in the sample is approximately distributed as Binomial($n = 100, p = 30/10,000 = 0.003$).

If $\frac{n}{N} \leq 0.05$, can use Binomial to approximate

Let, $Y \sim \text{HYG}(N, M, n)$
if, $\frac{n}{N} \leq 0.05$

then,

$Y \sim \text{Bin}(n, M/N)$

or, $Y \sim \text{Bin}(n, p)$
where,
 $p = M/N$

When your sample size is less than or equal to 5% of your population, the change in probabilities at each trial from not replacing is negligible \Rightarrow can treat the trials as independent.

Hypergeometric($N=100, M=30, n=10$) pmf

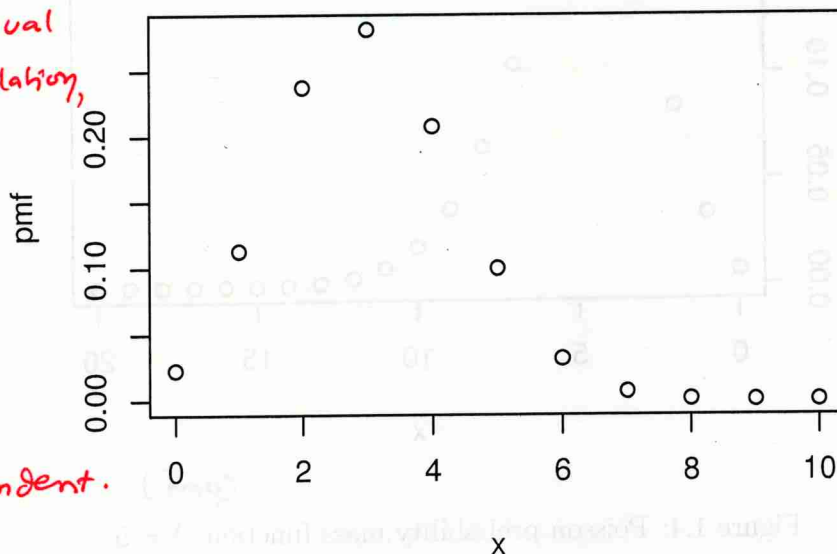


Figure 1.3: Hypergeometric probability mass function, $N=100, M=30, n=10$

6. Poisson (λ); $f(x) = \frac{\lambda^x e^{-\lambda}}{x!} \quad x = 0, 1, \dots$

Discrete

Models the number of occurrences of random events in space and time, where the average rate, λ per unit time (or area, or volume) is constant.

Recall \leftarrow

$\text{pois}(np) \approx \text{Bin}(n, p)$

When $n \gg 20$ and

$np \leq 5$

large n ,
small p

Let $X = \#$ events in t units of time
Then $X \sim \text{Poisson}(\lambda t)$

Example: Let X = the number of customers arriving at a bank in a given one hour time interval.

What is the sample space?

$E(X) = \lambda$, $\text{Var}(X) = \lambda$ $\Omega = \{\text{Nobody, one person, } \dots\}$

Poisson($\lambda=5$) pmf

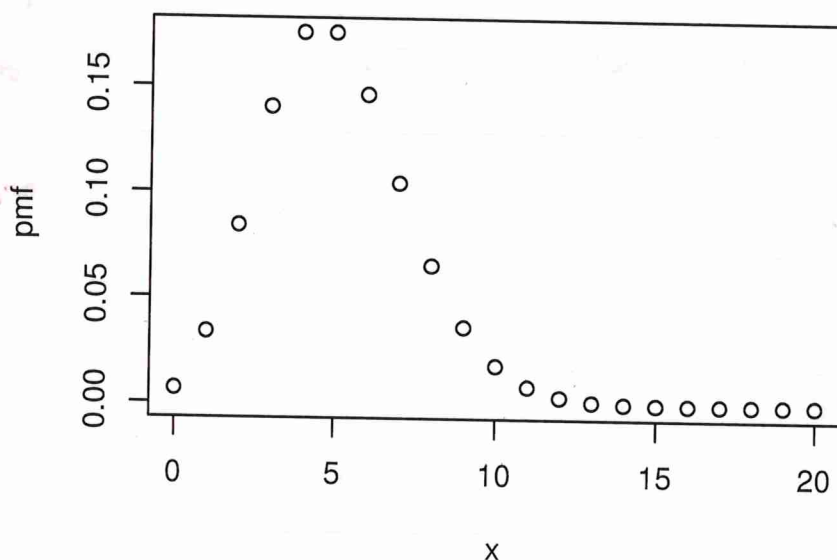


Figure 1.4: Poisson probability mass function, $\lambda = 5$

7. **Exponential** (mean θ); $f(x) = \frac{1}{\theta} e^{-x/\theta}$ $x > 0$ $\theta > 0$

Continuous!

Models lifetimes where there is no deterioration with age - or - waiting times between successive random events in a Poisson process. We also parameterize the exponential distribution using the rate parameter, $\lambda = 1/\theta$.

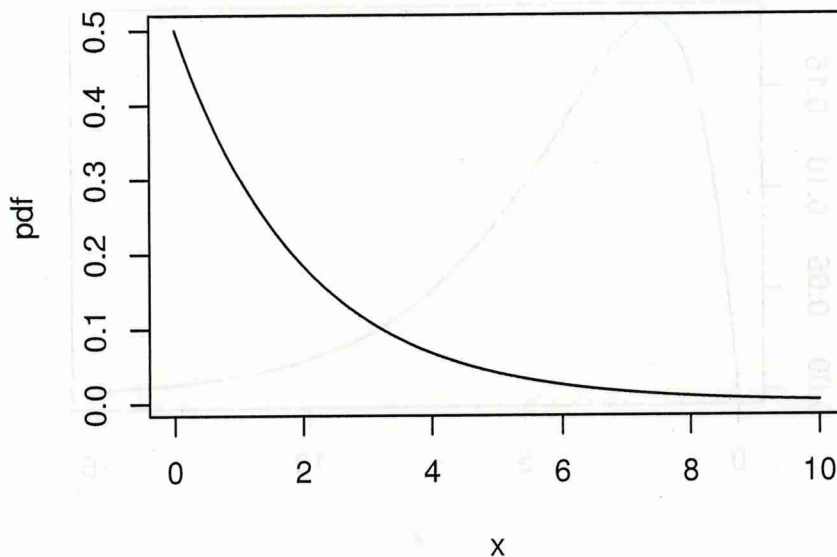
$$E(X) = \theta, \text{Var}(X) = \theta^2$$

Memoryless property \leftarrow

$$P(X > a+b | X > a) = P(X > b), \quad a, b > 0$$

If, $X \sim \text{pois}(\alpha)$ is counting the # of customers entering in an hour, then $Y \sim \text{Exp}(\text{mean } 1/\alpha)$ denotes the time between consecutive entrances.

Exponential(rate=.5) prob density function

Figure 1.5: Exponential density function, $\theta = 1/.5 = 2$

$$\lambda = \text{rate} = .5$$

8. Gamma(α, β); $f(x) = \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} \exp(-x/\beta)$, $x > 0$, $\alpha, \beta > 0$.

$$E(X) = \alpha\beta, \text{Var}(X) = \alpha\beta^2$$

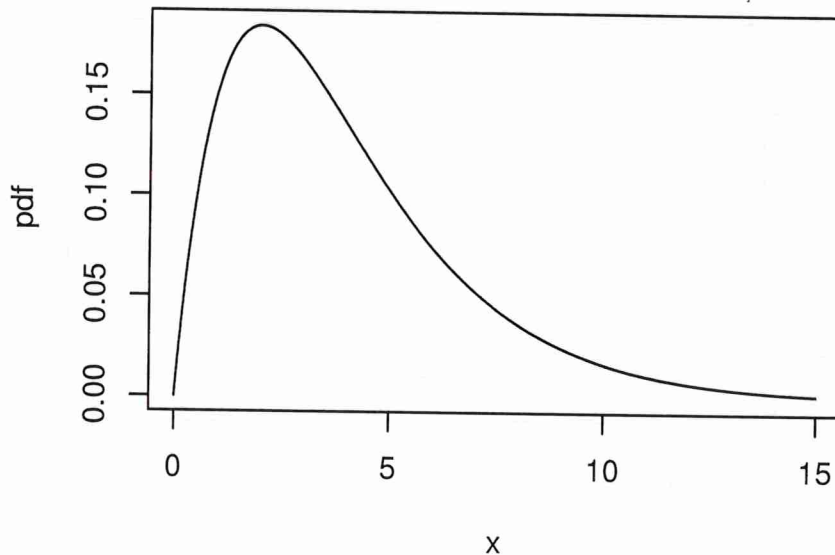
flexible model used to model lifetimes.

When $\alpha=1$, what does this simplify to?

$$\begin{aligned} \text{Gamma}(1, \beta) &\Rightarrow f(x) = \frac{1}{\Gamma(1)\beta^1} x^{1-1} \exp(-x/\beta) \\ &= \frac{1}{\beta} e^{-x/\beta} \end{aligned}$$

\therefore This becomes, $\exp(\text{mean } \beta)$
or, $\exp(\text{rate } 1/\beta)$

Gamma(alpha=2, beta=2) prob density function

Figure 1.6: Gamma density function, $\alpha = 2, \beta = 2$

9. Normal (μ, σ^2) ; $f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right\}$, $x, \mu \in \mathbb{R}$, $\sigma^2 > 0$.

Many measurements are approximately normal.

If $X \sim N(\mu, \sigma^2)$,

$$Z = \frac{X - \mu}{\sigma} \sim N(0, 1).$$

(Standard Normal Distribution)

mean (μ) ("mu")

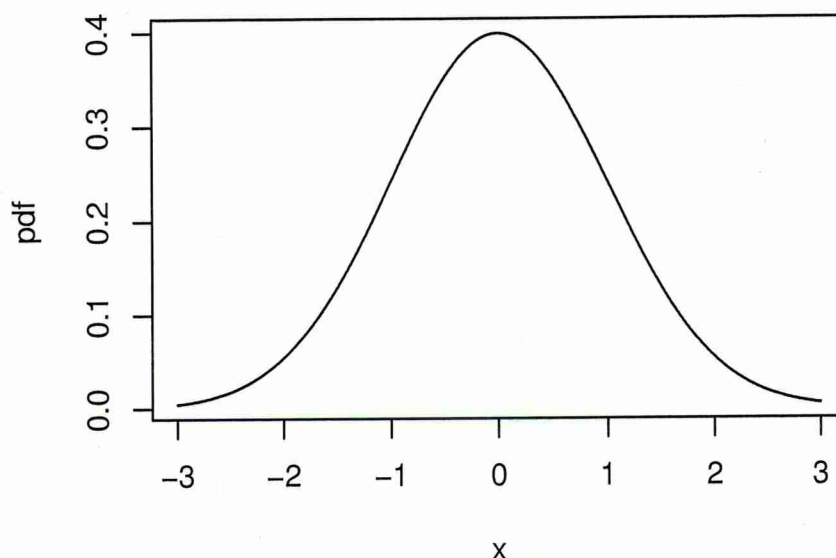
Variance (σ^2) ("sigma")

$$E(X) = \mu, \quad \text{Var}(X) = \sigma^2, \quad \text{sd} = \sigma = \text{standard deviation}$$

We assume $\sigma^2 > 0$, because $\sigma^2 = 0$ is trivial

\hookrightarrow leads to degenerated dist'n.

Normal(mean=0, sd=1) prob density function

Figure 1.7: Normal density function, $\mu = 0, \sigma = 1$

① If, $X_1 \sim N(\mu_1, \sigma_1^2)$
 $X_2 \sim N(\mu_2, \sigma_2^2)$

Notation: $X_1 \perp X_2$

Note: If X_1, \dots, X_n are independent with $X_i \sim N(\mu_i, \sigma_i^2)$ and a_1, \dots, a_n are constants,

$X_1 + X_2 \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$
 then $\sum_{i=1}^n a_i X_i \sim N(\sum a_i \mu_i, \sum a_i^2 \sigma_i^2)$

★ Linear combinations of Normal R.V's are Normal.

$X_1 + 2X_2 \sim N(\mu_1 + 2\mu_2, \sigma_1^2 + 4\sigma_2^2)$

Central Limit Theorem (CLT)

Let $S_n = \sum_{i=1}^n X_i$ be the sum of n independent random variables each with mean μ , variance σ^2 . Then

$$\frac{S_n - n\mu}{\sigma\sqrt{n}} \approx N(0, 1) \text{ for large } n,$$

where \approx means approximately distributed as.

Alternatively, the CLT is expressed as
 finite variance $\sigma^2 < \infty$.
 $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \approx N(0, 1)$

From Jan 9 lecture [maybe helpful for your
Lecture Assignment-1]

Let, X be a R.V (random variable)

1) If, X is discrete, the distribution function is called probability mass function (pmf)

$$f(x) = P(X=x) ; \forall x \in \mathcal{X}$$

for all x belonging to \mathcal{X}

$f(x)$ has these constraints

① $f(x) \geq 0 ; \forall x \in \mathcal{X}$

② $\sum_{x \in \mathcal{X}} f(x) = 1.$

2) If, X is Continuous, then it can take on values in an interval {subset of \mathbb{R} } and we call the distribution function 'probability density function' (pdf)

Consider ; $F(x) = P(X \leq x) ; \forall x \in \mathcal{X}$

The pdf is defined as : $f(x) = \frac{d}{dx} F(x)$

\Leftrightarrow

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

Fundamental
Theorem
of
Calculus

$f(x)$ has the following constraints :

① $f(x) \geq 0 ; \forall x \in \mathcal{X}$

② $\int_{-\infty}^{\infty} f(x) dx = 1$

③ $P(X=x) = 0 ; \forall x \in \mathcal{X}$

What does it mean for $X=x$?

\hookrightarrow Both $X \geq x$ and $X \leq x$ are true

$$\begin{aligned} \text{So, } P(X=x) &= P(x \leq X \leq x) \\ &= F(x) - F(x) = 0 \end{aligned}$$

Consider a bin of $N=100$ marbles, with $M=30$ white marbles and $N-M=70$ black marbles. Suppose, we sample $n=40$ marbles from the bin without replacement. How many possible outcomes make up the sample space of this experiment? How about with $n=10$ marbles?

Sol:- ~~For~~ ^{For} $\boxed{n=40}$

2 types of marbles $\begin{cases} \rightarrow \text{black (x 70) B} \\ \rightarrow \text{white (x 30) W} \end{cases}$

A string of length 40 consisting of B's and W's

or, may be no W's, only B's, because $\# \text{ white} \leq n$ ($n=40$)

Let, x be the # white marbles drawn

$$\therefore 0 \leq x \leq 30$$

For each x , there are $\binom{n}{x}$ ways to arrange the white and black marbles,

$$\text{Therefore, } |\Omega| = \sum_{x=0}^{30} \binom{n}{x} = \sum_{x=0}^{30} \binom{40}{x} = \sum_{x=0}^{30} \frac{40!}{x!(40-x)!}$$

$$\text{i.e., Case 1; } x=0, \quad \frac{40!}{0!(40-0)!} = 1$$

$$\text{Case 2; } x=1, \quad \frac{40!}{1!39!} = 40$$

\vdots

$$\text{Case 31; } x=30, \quad \frac{40!}{30!10!}$$

Add up all the cases

$$= \sum_{x=0}^{30} \frac{40!}{x!(40-x)!} \quad (\text{Total Outcomes})$$

should be.

$$\text{For, } \boxed{n=10} \quad (x \leq n)$$

$$\text{Similarly, } \sum_{x=0}^{10} \frac{10!}{x!(10-x)!} = 1024 = 2^{10} \quad (\text{Total Outcomes})$$

In this case, we can simply think, both the # white and Black Marbles are $> n$, so, there will be 2^n outcomes.

Let see, how this problem is ^{related to} a hypergeometric distribution,

What is the probability of getting 6 white balls

a) in the sample of size 10 ?

b) in the sample of size 40 ?

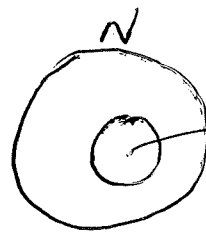
↳ Here, $N = 100$, $M = 30$ (white balls)

a) $n = 10$; b) $n = 40$

$$P(X=6) = \frac{\binom{30}{6} \binom{70}{4}}{\binom{100}{10}} ; P(X=6) = \frac{\binom{30}{6} \binom{70}{36}}{\binom{100}{40}}$$

Hypergeometric
Distribution

$N \begin{cases} M \\ N-M \end{cases}$



sample of size $n \begin{cases} x_s \\ n - x_s \end{cases}$

⑥ "N" contains "M" number of Success "S"

then, there are " $N-M$ " number of Failure "F"

we are taking a sample of size " n " from " N "

and, we are interested in how many Success "S"
are there in the sample of size "n"

∴ let X be a r.v, which defines the # S
in the sample of size n .

$$\therefore P(X=x) = \frac{\binom{M}{x} \binom{N-M}{n-x}}{\binom{N}{n}}$$