

# **Statistical and Mathematical Methods for Data Analysis**

**Dr. Syed Faisal Bukhari**

Associate Professor

Department of Data Science

Faculty of Computing and Information Technology

University of the Punjab

# Textbooks

- ❑ **Probability & Statistics for Engineers & Scientists**, Ninth Edition, Ronald E. Walpole, Raymond H. Myer
- ❑ **Elementary Statistics: Picturing the World**, 6<sup>th</sup> Edition, Ron Larson and Betsy Farber
- ❑ **Elementary Statistics**, 13<sup>th</sup> Edition, Mario F. Triola

# Reference books

- ❑ **Probability Demystified**, Allan G. Bluman
- ❑ **Schaum's Outline of Probability and Statistics**
- ❑ **MATLAB Primer**, Seventh Edition
- ❑ **MATLAB Demystified** by McMahan, David

# Reference books

- ❑ **Probability and Statistical Inference, Ninth Edition,** Robert V. Hogg, Elliot A. Tanis, Dale L. Zimmerman
- ❑ **Probability Demystified,** Allan G. Bluman
- ❑ **Practical Statistics for Data Scientists: 50 Essential Concepts,** Peter Bruce and Andrew Bruce
- ❑ **Schaum's Outline of Probability,** Second Edition, Seymour Lipschutz, Marc Lipson
- ❑ **Python for Probability, Statistics, and Machine Learning,** José Unpingco

# References

Readings for these lecture notes:

❑ **Probability & Statistics for Engineers & Scientists**,  
Ninth edition, Ronald E. Walpole, Raymond H.  
Myer

❑ **Elementary Statistics**, 10<sup>th</sup> Edition, Mario F. Triola

❑ **Probability Demystified**, Allan G. Bluman

These notes contain material from the above three books.

**“A goal is a dream with a deadline.”**

**— Napoleon Hill**

# Hypergeometric Distribution

- ❑ **Hypergeometric Distribution** If we sample from a small finite population **without replacement**, the binomial distribution should not be used because the events are **not independent**.
- ❑ If sampling is done **without replacement** and the outcomes belong to **one of two types**, we can use the **hypergeometric distribution**

# Hypergeometric Distribution

- ❑ The simplest way to view the **distinction** between the binomial distribution and the hypergeometric distribution is to note the **way the sampling is done**.
- ❑ The types of applications for the **hypergeometric** are very similar to those for the binomial distribution. We are interested in computing probabilities for the number of observations that **fall into a particular category**.



# Applications

- ❑ Applications for the hypergeometric distribution are found in many areas, with heavy use in **acceptance sampling, electronic testing, and quality assurance**. Obviously, in many of these fields, **testing is done at the expense of the item being tested**.
- ❑ That is, the item is **destroyed** and hence **cannot be replaced** in the sample

# Hypergeometric Distribution [1]

A **hypergeometric experiment** has the following properties:

1. Each trial of an experiment results in **an outcome** that can be classified into one of the two categories **success or failure**.
2. The successive trials are **dependent**.
3. The probability of success **changes** from trial to trial.
4. The experiment is repeated **a fixed number** of times.

# Hypergeometric Distribution

- In general, we are interested in the probability of selecting  $x$  successes from the  $k$  items labeled successes and  $n - x$  failures from the  $N - k$  items labeled failures when a random sample of size  $n$  is selected from  $N$  items.
- This is known as a **hypergeometric experiment**, that is, one that possesses the following two properties:
  1. A random sample of size  $n$  is selected **without replacement** from  $N$  items.
  2. Of the  $N$  items,  $k$  may be classified as **successes** and  $N - k$  are classified as **failures**.

# Hypergeometric Distribution

- ❑ The number  **$X$  of successes** of a hypergeometric experiment is called a **hypergeometric random variable**.
- ❑ Accordingly, the probability distribution of the hypergeometric variable is called the **hypergeometric distribution**, and its values are denoted by  **$h(x; N, n, k)$** , since they depend on the **number of successes  $k$**  in the set  **$N$**  from which we **select  $n$  items**.

# Hypergeometric Distribution [2]

This distribution is the case of sampling **without replacement**. The formula to calculate probabilities is given by

$$P(X = x) = h(x; N, n, k) \\ = \binom{k}{x} \binom{N-k}{n-x} / \binom{N}{n}, \\ \max\{0, n - (N-k)\} \leq x \leq \min\{n, k\}$$

OR

$$h(x; N, n, k) = \frac{\binom{k}{x} \binom{N-k}{n-x}}{\binom{N}{n}}, \max\{0, n - (N-k)\} \leq x \leq \min\{n, k\}$$

# Hypergeometric Distribution [3]

- ❑ It has **three** parameters i.e.,  $N$ ,  $n$ , and  $k$
- ❑  **$N$** : The number of items in the **population**
- ❑  **$k$** : The number of items in the **population** that are classified as **successes**.
- ❑  **$n$** : The number of items in the sample
- ❑  **$x$** : The number of items in the **sample** that are classified as **successes**.

# Hypergeometric Distribution [4]

**Example1:** Suppose we randomly select **5** cards **without replacement** from an ordinary deck of playing cards. What is the probability of getting exactly **2 red cards**? Also, implement it in Python.

**Solution:** This is a hypergeometric experiment in which we know the following:

**N = 52;** since there are 52 cards in a deck.

**k = 26;** since there are 26 red cards in a deck.

**n = 5;** since we randomly select 5 cards from the deck.

**x = 2;** since 2 of the cards we select are red.

$$h(x; N, n, k) = \frac{{}_k C_x {}_{N-k} C_{n-x}}{{}_N C_n}$$

$$h(2; 52, 5, 26) = \frac{{}_{26} C_2 {}_{26} C_3}{{}_{52} C_5}$$

$$h(2; 52, 5, 26) = (325) (2600) / (2,598,960) = 0.325 \text{ or } 32.51\%$$



## Python code:

```
from scipy.stats import hypergeom
```

```
N = 52      # Total number of cards
```

```
n = 5       # Total number of cards randomly selected
```

```
k = 26      # Since there are 26 red cards
```

```
x = 2       # Two red cards
```

```
prob= round(hypergeom.pmf(x, N, n, k), 4)
```

```
# Compute probabilities corresponding to random variable x
```

```
print('Probability is :', prob)
```

```
# Probability is : 0.3251
```

# Hypergeometric Distribution [5]

**Example:** A committee of 4 people is selected at random without replacement from a group of 6 men and 4 women. Find the probability that the committee consists of 2 men and 2 women.

## Solution:

$$P(X = x) = h(x; N, n, k) = \frac{{}_k C_x {}_{N-k} C_{n-x}}{{}_N C_n}, \max\{0, n - (N - k)\} \leq x \leq \min\{n, k\}$$

**$N = 10; k = 6; n = 4; x = 2$  (let  $x$  denotes number of men)**

$$h(x; N, n, k) = \frac{{}_6 C_x {}_4 C_{4-x}}{{}_{10} C_4}$$

$$h(2; 10, 4, 6) = \frac{{}_6 C_2 {}_4 C_2}{{}_{10} C_4}$$

$$h(2; 10, 4, 6) = \frac{(15)(6)}{(210)} = 0.429 \text{ or } 42.9\%$$

# Hypergeometric Distribution [6]

**Example:** A lot of **12 oxygen tanks** contains **3** defective ones. If **4 tanks** are randomly selected and tested, find the probability that exactly **one will be defective**.

## Solution:

$$P(X = x) = h(x; N, n, k) = \frac{{}_k C_x {}_{N-k} C_{n-x}}{{}_N C_n}, \max\{0, n - (N - k)\} \leq x \leq \min\{n, k\}$$

OR

$$h(x; N, n, k) = \frac{{}_k C_x {}_{N-k} C_{n-x}}{{}_N C_n}, \max\{0, n - (N - k)\} \leq x \leq \min\{n, k\}$$

**$N = 12; k = 3; n = 4; x = 1$  (let  $x$  denotes defective tanks)**

$$P(X = 1) = \frac{{}_3 C_1 {}_9 C_3}{{}_{12} C_4}$$

$$P(X = 1) = \frac{(3)(84)}{(495)} = 0.509 \text{ or } 50.9\%$$

# Hypergeometric Distribution [1]

**Example:** In a box of **12** shirts there are **5** defective ones. If **5** shirts are sold at random, find the probability that exactly two are defective.

## Solution:

$$h(x; N, n, k) = \frac{\binom{k}{x} \binom{N-k}{n-x}}{\binom{N}{n}}, \max\{0, n-(N-k)\} \leq x \leq \min\{n, k\}$$

Defective shirts	Non-detective shirts	Total
5	7	12

$N = 12, k = 5, n = 5, \text{ and } x = 2$

Let **X** denotes the number of **defective shirts**

$$P(X = 2) = \binom{5}{2} \binom{7}{3} / \binom{12}{5} = 0.442$$

# Hypergeometric Distribution [2]

**Example:** In a fitness club of **18** members, **10** prefer the exercise bicycle and **8** prefer the aerobic stepper. If **6** members are selected at random, find the probability that exactly **3** use the bicycle.



## Solution:

$$h(x; N, n, k) = \frac{{}_k C_x {}_{N-k} C_{n-x}}{{}_N C_n}, \max\{0, n-(N-k)\} \leq x \leq \min\{n, k\}$$

Exercise Bicycle	Aerobic Stepper	Total
10	8	18

$N = 18, k = 10, n = 6, \text{ and } x = 3$

Let **X** denotes the number of **bicycles**

$$P(X=3) = ({}_{10} C_3)({}_8 C_3) / {}_{18} C_6 = 0.362$$

# Hypergeometric Distribution [3]

**Example:** In a shipment of **10** lawn chairs, **6** are brown and **4** are blue. If **3** chairs are sold at random, find the probability that all are **brown**.

## Solution:

$$h(x; N, n, k) = \frac{{}_k C_x {}_{N-k} C_{n-x}}{{}_N C_n}, \max\{0, n-(N-k)\} \leq x \leq \min\{n, k\}$$

Brown	Blue	Total
6	4	10

$N = 10$ ,  $k = 6$ ,  $n = 3$ , and  $x = 3$

Let  $X$  denotes the number of **brown** chairs

$$P(X = 3) = {}_6 C_3 {}_4 C_0 / {}_{10} C_3 = 0.167$$

# Hypergeometric Distribution [4]

**Example:** A class consists of 5 women and 4 men. If a committee of 3 people is selected at random without replacement, find the probability that all 3 are women.

## Solution:

$$h(x; N, n, k) = \frac{{}_k C_x {}_{N-k} C_{n-x}}{{}_N C_n}, \max\{0, n-(N-k)\} \leq x \leq \min\{n, k\}$$

Men	Women	Total
4	5	9

$N = 9, k = 5, n = 3$ , and  $x = 3$

Let **X** denotes the number of **women**

$$P(X=3) = ({}_5 C_3)({}_4 C_0) / {}_9 C_3 = 0.119$$

# Hypergeometric Distribution [5]

**Example:** A box contains **3 red** balls and **3 white balls**. If **two balls** are selected at random without replacement, find the probability that both are **red**.

## Solution:

$$h(x; N, n, k) = \frac{{}_k C_x {}_{N-k} C_{n-x}}{{}_N C_n}, \max\{0, n-(N-k)\} \leq x \leq \min\{n, k\}$$

Red	White	Total
3	3	6

$N = 6, k = 3, n = 2$ , and  $x = 2$

Let **X** denotes the number of **red balls**

$$P(X = 2) = {}_3 C_2 {}_3 C_0 / {}_6 C_2 = 0.2$$

Table A.1 Binomial Probability Sums  $\sum_{x=0}^r b(x; n, p)$

<i>n</i>	<i>r</i>	<i>p</i>									
		0.10	0.20	0.25	0.30	0.40	0.50	0.60	0.70	0.80	0.90
<b>1</b>	<b>0</b>	0.9000	0.8000	0.7500	0.7000	0.6000	0.5000	0.4000	0.3000	0.2000	0.1000
	<b>1</b>	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
<b>2</b>	<b>0</b>	0.8100	0.6400	0.5625	0.4900	0.3600	0.2500	0.1600	0.0900	0.0400	0.0100
	<b>1</b>	0.9900	0.9600	0.9375	0.9100	0.8400	0.7500	0.6400	0.5100	0.3600	0.1900
	<b>2</b>	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
<b>3</b>	<b>0</b>	0.7290	0.5120	0.4219	0.3430	0.2160	0.1250	0.0640	0.0270	0.0080	0.0010
	<b>1</b>	0.9720	0.8960	0.8438	0.7840	0.6480	0.5000	0.3520	0.2160	0.1040	0.0280
	<b>2</b>	0.9990	0.9920	0.9844	0.9730	0.9360	0.8750	0.7840	0.6570	0.4880	0.2710
	<b>3</b>	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
<b>4</b>	<b>0</b>	0.6561	0.4096	0.3164	0.2401	0.1296	0.0625	0.0256	0.0081	0.0016	0.0001
	<b>1</b>	0.9477	0.8192	0.7383	0.6517	0.4752	0.3125	0.1792	0.0837	0.0272	0.0037
	<b>2</b>	0.9963	0.9728	0.9492	0.9163	0.8208	0.6875	0.5248	0.3483	0.1808	0.0523
	<b>3</b>	0.9999	0.9984	0.9961	0.9919	0.9744	0.9375	0.8704	0.7599	0.5904	0.3439
	<b>4</b>	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
<b>5</b>	<b>0</b>	0.5905	0.3277	0.2373	0.1681	0.0778	0.0313	0.0102	0.0024	0.0003	0.0000
	<b>1</b>	0.9185	0.7373	0.6328	0.5282	0.3370	0.1875	0.0870	0.0308	0.0067	0.0005
	<b>2</b>	0.9914	0.9421	0.8965	0.8369	0.6826	0.5000	0.3174	0.1631	0.0579	0.0086
	<b>3</b>	0.9995	0.9933	0.9844	0.9692	0.9130	0.8125	0.6630	0.4718	0.2627	0.0815
	<b>4</b>	1.0000	0.9997	0.9990	0.9976	0.9898	0.9688	0.9222	0.8319	0.6723	0.4095
	<b>5</b>	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
<b>6</b>	<b>0</b>	0.5314	0.2621	0.1780	0.1176	0.0467	0.0156	0.0041	0.0007	0.0001	0.0000
	<b>1</b>	0.8857	0.6554	0.5339	0.4202	0.2333	0.1094	0.0410	0.0109	0.0016	0.0001
	<b>2</b>	0.9842	0.9011	0.8306	0.7443	0.5443	0.3438	0.1792	0.0705	0.0170	0.0013
	<b>3</b>	0.9987	0.9830	0.9624	0.9295	0.8208	0.6563	0.4557	0.2557	0.0989	0.0159
	<b>4</b>	0.9999	0.9984	0.9954	0.9891	0.9590	0.8906	0.7667	0.5798	0.3446	0.1143
	<b>5</b>	1.0000	0.9999	0.9998	0.9993	0.9959	0.9844	0.9533	0.8824	0.7379	0.4686
	<b>6</b>	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
<b>7</b>	<b>0</b>	0.4783	0.2097	0.1335	0.0824	0.0280	0.0078	0.0016	0.0002	0.0000	
	<b>1</b>	0.8503	0.5767	0.4449	0.3294	0.1586	0.0625	0.0188	0.0038	0.0004	0.0000
	<b>2</b>	0.9743	0.8520	0.7564	0.6471	0.4199	0.2266	0.0963	0.0288	0.0047	0.0002
	<b>3</b>	0.9973	0.9667	0.9294	0.8740	0.7102	0.5000	0.2898	0.1260	0.0333	0.0027
	<b>4</b>	0.9996	0.9953	0.9871	0.9712	0.9037	0.7734	0.5801	0.3529	0.1480	0.0257
	<b>5</b>	1.0000	0.9996	0.9967	0.9962	0.9812	0.9375	0.8414	0.6706	0.4233	0.1497
	<b>6</b>		1.0000	0.9999	0.9998	0.9984	0.9922	0.9720	0.9176	0.7903	0.5217
	<b>7</b>			1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000



Table A.1 (continued) Binomial Probability Sums  $\sum_{x=0}^r b(x; n, p)$ 

n	r	p									
		0.10	0.20	0.25	0.30	0.40	0.50	0.60	0.70	0.80	0.90
8	0	0.4305	0.1678	0.1001	0.0576	0.0168	0.0039	0.0007	0.0001	0.0000	
	1	0.8131	0.5033	0.3671	0.2553	0.1064	0.0352	0.0085	0.0013	0.0001	
	2	0.9619	0.7969	0.6785	0.5518	0.3154	0.1445	0.0498	0.0113	0.0012	0.0000
	3	0.9950	0.9437	0.8862	0.8059	0.5941	0.3633	0.1737	0.0580	0.0104	0.0004
	4	0.9996	0.9896	0.9727	0.9420	0.8263	0.6367	0.4059	0.1941	0.0563	0.0050
	5	1.0000	0.9988	0.9958	0.9887	0.9502	0.8555	0.6846	0.4482	0.2031	0.0381
	6		0.9999	0.9996	0.9987	0.9915	0.9648	0.8936	0.7447	0.4967	0.1869
	7		1.0000	1.0000	0.9999	0.9993	0.9961	0.9832	0.9424	0.8322	0.5695
	8				1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
9	0	0.3874	0.1342	0.0751	0.0404	0.0101	0.0020	0.0003	0.0000		
	1	0.7748	0.4362	0.3003	0.1980	0.0705	0.0195	0.0038	0.0004	0.0000	
	2	0.9470	0.7382	0.6007	0.4628	0.2318	0.0898	0.0250	0.0043	0.0003	0.0000
	3	0.9917	0.9144	0.8343	0.7297	0.4826	0.2539	0.0994	0.0253	0.0031	0.0001
	4	0.9991	0.9804	0.9511	0.9012	0.7334	0.5000	0.2666	0.0988	0.0196	0.0009
	5	0.9999	0.9969	0.9900	0.9747	0.9006	0.7461	0.5174	0.2703	0.0856	0.0083
	6	1.0000	0.9997	0.9987	0.9957	0.9750	0.9102	0.7682	0.5372	0.2618	0.0530
	7		1.0000	0.9999	0.9996	0.9962	0.9805	0.9295	0.8040	0.5638	0.2252
	8			1.0000	1.0000	0.9997	0.9980	0.9899	0.9596	0.8658	0.6126
	9					1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
10	0	0.3487	0.1074	0.0563	0.0282	0.0060	0.0010	0.0001	0.0000		
	1	0.7361	0.3758	0.2440	0.1493	0.0464	0.0107	0.0017	0.0001	0.0000	
	2	0.9298	0.6778	0.5256	0.3828	0.1673	0.0547	0.0123	0.0016	0.0001	
	3	0.9872	0.8791	0.7759	0.6496	0.3823	0.1719	0.0548	0.0106	0.0009	0.0000
	4	0.9984	0.9672	0.9219	0.8497	0.6331	0.3770	0.1662	0.0473	0.0064	0.0001
	5	0.9999	0.9936	0.9803	0.9527	0.8338	0.6230	0.3669	0.1503	0.0328	0.0016
	6	1.0000	0.9991	0.9965	0.9894	0.9452	0.8281	0.6177	0.3504	0.1209	0.0128
	7		0.9999	0.9996	0.9984	0.9877	0.9453	0.8327	0.6172	0.3222	0.0702
	8		1.0000	1.0000	0.9999	0.9983	0.9893	0.9536	0.8507	0.6242	0.2639
	9				1.0000	0.9999	0.9990	0.9940	0.9718	0.8926	0.6513
	10					1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
11	0	0.3138	0.0859	0.0422	0.0198	0.0036	0.0005	0.0000			
	1	0.6974	0.3221	0.1971	0.1130	0.0302	0.0059	0.0007	0.0000		
	2	0.9104	0.6174	0.4552	0.3127	0.1189	0.0327	0.0059	0.0006	0.0000	
	3	0.9815	0.8389	0.7133	0.5696	0.2963	0.1133	0.0293	0.0043	0.0002	
	4	0.9972	0.9496	0.8854	0.7897	0.5328	0.2744	0.0994	0.0216	0.0020	0.0000
	5	0.9997	0.9883	0.9657	0.9218	0.7535	0.5000	0.2465	0.0782	0.0117	0.0003
	6	1.0000	0.9980	0.9924	0.9784	0.9006	0.7256	0.4672	0.2103	0.0504	0.0028
	7		0.9998	0.9988	0.9957	0.9707	0.8867	0.7037	0.4304	0.1611	0.0185
	8		1.0000	0.9999	0.9994	0.9941	0.9673	0.8811	0.6873	0.3826	0.0896
	9			1.0000	1.0000	0.9993	0.9941	0.9698	0.8870	0.6779	0.3026
	10					1.0000	0.9995	0.9964	0.9802	0.9141	0.6862
	11						1.0000	1.0000	1.0000	1.0000	1.0000

Table A.1 (continued) Binomial Probability Sums  $\sum_{x=0}^r b(x; n, p)$ 

n	r	p									
		0.10	0.20	0.25	0.30	0.40	0.50	0.60	0.70	0.80	0.90
12	0	0.2824	0.0687	0.0317	0.0138	0.0022	0.0002	0.0000			
	1	0.6590	0.2749	0.1584	0.0850	0.0196	0.0032	0.0003	0.0000		
	2	0.8891	0.5583	0.3907	0.2528	0.0834	0.0193	0.0028	0.0002	0.0000	
	3	0.9744	0.7946	0.6488	0.4925	0.2253	0.0730	0.0153	0.0017	0.0001	
	4	0.9957	0.9274	0.8424	0.7237	0.4382	0.1938	0.0573	0.0095	0.0006	0.0000
	5	0.9995	0.9806	0.9456	0.8822	0.6652	0.3872	0.1582	0.0386	0.0039	0.0001
	6	0.9999	0.9961	0.9857	0.9614	0.8418	0.6128	0.3348	0.1178	0.0194	0.0005
	7	1.0000	0.9994	0.9972	0.9905	0.9427	0.8062	0.5618	0.2763	0.0726	0.0043
	8		0.9999	0.9996	0.9983	0.9847	0.9270	0.7747	0.5075	0.2054	0.0256
	9		1.0000	1.0000	0.9998	0.9972	0.9807	0.9166	0.7472	0.4417	0.1109
	10				1.0000	0.9997	0.9968	0.9804	0.9150	0.7251	0.3410
	11					1.0000	0.9998	0.9978	0.9862	0.9313	0.7176
	12						1.0000	1.0000	1.0000	1.0000	1.0000
13	0	0.2542	0.0550	0.0238	0.0097	0.0013	0.0001	0.0000			
	1	0.6213	0.2336	0.1267	0.0637	0.0126	0.0017	0.0001	0.0000		
	2	0.8661	0.5017	0.3326	0.2025	0.0579	0.0112	0.0013	0.0001		
	3	0.9658	0.7473	0.5843	0.4206	0.1686	0.0461	0.0078	0.0007	0.0000	
	4	0.9935	0.9009	0.7940	0.6543	0.3530	0.1334	0.0321	0.0040	0.0002	
	5	0.9991	0.9700	0.9198	0.8346	0.5744	0.2905	0.0977	0.0182	0.0012	0.0000
	6	0.9999	0.9930	0.9757	0.9376	0.7712	0.5000	0.2288	0.0624	0.0070	0.0001
	7	1.0000	0.9988	0.9944	0.9818	0.9023	0.7095	0.4256	0.1654	0.0300	0.0009
	8		0.9998	0.9990	0.9960	0.9679	0.8666	0.6470	0.3457	0.0991	0.0065
	9		1.0000	0.9999	0.9993	0.9922	0.9539	0.8314	0.5794	0.2527	0.0342
	10			1.0000	0.9999	0.9987	0.9888	0.9421	0.7975	0.4983	0.1339
	11				1.0000	0.9999	0.9983	0.9874	0.9363	0.7664	0.3787
	12					1.0000	0.9999	0.9987	0.9903	0.9450	0.7458
	13						1.0000	1.0000	1.0000	1.0000	1.0000
14	0	0.2288	0.0440	0.0178	0.0068	0.0008	0.0001	0.0000			
	1	0.5846	0.1979	0.1010	0.0475	0.0081	0.0009	0.0001			
	2	0.8416	0.4481	0.2811	0.1608	0.0398	0.0065	0.0006	0.0000		
	3	0.9559	0.6982	0.5213	0.3552	0.1243	0.0287	0.0039	0.0002		
	4	0.9908	0.8702	0.7415	0.5842	0.2793	0.0898	0.0175	0.0017	0.0000	
	5	0.9985	0.9561	0.8883	0.7805	0.4859	0.2120	0.0583	0.0083	0.0004	
	6	0.9998	0.9884	0.9617	0.9067	0.6925	0.3953	0.1501	0.0315	0.0024	0.0000
	7	1.0000	0.9976	0.9897	0.9685	0.8499	0.6047	0.3075	0.0933	0.0116	0.0002
	8		0.9996	0.9978	0.9917	0.9417	0.7880	0.5141	0.2195	0.0439	0.0015
	9		1.0000	0.9997	0.9983	0.9825	0.9102	0.7207	0.4158	0.1298	0.0092
	10			1.0000	0.9998	0.9961	0.9713	0.8757	0.6448	0.3018	0.0441
	11				1.0000	0.9994	0.9935	0.9602	0.8392	0.5519	0.1584
	12					0.9999	0.9991	0.9919	0.9525	0.8021	0.4154
	13					1.0000	0.9999	0.9992	0.9932	0.9560	0.7712
	14						1.0000	1.0000	1.0000	1.0000	1.0000

Table A.1 (continued) Binomial Probability Sums  $\sum_{x=0}^r b(x; n, p)$

n	r	p									
		0.10	0.20	0.25	0.30	0.40	0.50	0.60	0.70	0.80	0.90
15	0	0.2059	0.0352	0.0134	0.0047	0.0005	0.0000				
	1	0.5490	0.1671	0.0802	0.0353	0.0052	0.0005	0.0000			
	2	0.8159	0.3980	0.2361	0.1268	0.0271	0.0037	0.0003	0.0000		
	3	0.9444	0.6482	0.4613	0.2969	0.0905	0.0176	0.0019	0.0001		
	4	0.9873	0.8358	0.6865	0.5155	0.2173	0.0592	0.0093	0.0007	0.0000	
	5	0.9978	0.9389	0.8516	0.7216	0.4032	0.1509	0.0338	0.0037	0.0001	
	6	0.9997	0.9819	0.9434	0.8689	0.6098	0.3036	0.0950	0.0152	0.0008	
	7	1.0000	0.9958	0.9827	0.9500	0.7869	0.5000	0.2131	0.0500	0.0042	0.0000
	8		0.9992	0.9958	0.9848	0.9050	0.6904	0.3902	0.1311	0.0181	0.0003
	9		0.9999	0.9992	0.9963	0.9662	0.8491	0.5968	0.2784	0.0611	0.0022
	10		1.0000	0.9999	0.9993	0.9907	0.9408	0.7827	0.4845	0.1642	0.0127
	11			1.0000	0.9999	0.9981	0.9824	0.9095	0.7031	0.3518	0.0556
	12				1.0000	0.9997	0.9963	0.9729	0.8732	0.6020	0.1841
	13					1.0000	0.9995	0.9948	0.9647	0.8329	0.4510
	14						1.0000	0.9995	0.9953	0.9648	0.7941
	15							1.0000	1.0000	1.0000	1.0000
16	0	0.1853	0.0281	0.0100	0.0033	0.0003	0.0000				
	1	0.5147	0.1407	0.0635	0.0261	0.0033	0.0003	0.0000			
	2	0.7892	0.3518	0.1971	0.0994	0.0183	0.0021	0.0001			
	3	0.9316	0.5981	0.4050	0.2459	0.0651	0.0106	0.0009	0.0000		
	4	0.9830	0.7982	0.6302	0.4499	0.1666	0.0384	0.0049	0.0003		
	5	0.9967	0.9183	0.8103	0.6598	0.3288	0.1051	0.0191	0.0016	0.0000	
	6	0.9995	0.9733	0.9204	0.8247	0.5272	0.2272	0.0583	0.0071	0.0002	
	7	0.9999	0.9930	0.9729	0.9256	0.7161	0.4018	0.1423	0.0257	0.0015	0.0000
	8	1.0000	0.9985	0.9925	0.9743	0.8577	0.5982	0.2839	0.0744	0.0070	0.0001
	9		0.9998	0.9984	0.9929	0.9417	0.7728	0.4728	0.1753	0.0267	0.0005
	10		1.0000	0.9997	0.9984	0.9809	0.8949	0.6712	0.3402	0.0817	0.0033
	11			1.0000	0.9997	0.9951	0.9616	0.8334	0.5501	0.2018	0.0170
	12				1.0000	0.9991	0.9894	0.9349	0.7541	0.4019	0.0684
	13					0.9999	0.9979	0.9817	0.9006	0.6482	0.2108
	14					1.0000	0.9997	0.9967	0.9739	0.8593	0.4853
	15						1.0000	0.9997	0.9967	0.9719	0.8147
	16							1.0000	1.0000	1.0000	1.0000

Table A.1 (continued) Binomial Probability Sums  $\sum_{x=0}^r b(x; n, p)$

n	r	P									
		0.10	0.20	0.25	0.30	0.40	0.50	0.60	0.70	0.80	0.90
17	0	0.1088	0.0225	0.0075	0.0023	0.0002	0.0000				
	1	0.4818	0.1182	0.0501	0.0193	0.0021	0.0001	0.0000			
	2	0.7618	0.3096	0.1637	0.0774	0.0123	0.0012	0.0001			
	3	0.9174	0.5489	0.3530	0.2019	0.0464	0.0064	0.0005	0.0000		
	4	0.9779	0.7582	0.5739	0.3887	0.1260	0.0245	0.0025	0.0001		
	5	0.9953	0.8943	0.7653	0.5968	0.2639	0.0717	0.0106	0.0007	0.0000	
	6	0.9992	0.9623	0.8929	0.7752	0.4478	0.1662	0.0348	0.0032	0.0001	
	7	0.9999	0.9891	0.9598	0.8954	0.6405	0.3145	0.0919	0.0127	0.0005	
	8	1.0000	0.9974	0.9876	0.9597	0.8011	0.5000	0.1989	0.0403	0.0026	0.0000
	9		0.9995	0.9969	0.9873	0.9081	0.6855	0.3595	0.1046	0.0109	0.0001
	10		0.9999	0.9994	0.9968	0.9652	0.8338	0.5522	0.2248	0.0377	0.0008
	11		1.0000	0.9999	0.9993	0.9894	0.9283	0.7361	0.4032	0.1057	0.0047
	12			1.0000	0.9999	0.9975	0.9755	0.8740	0.6113	0.2418	0.0221
	13				1.0000	0.9995	0.9936	0.9536	0.7981	0.4511	0.0826
	14					0.9999	0.9988	0.9877	0.9226	0.6904	0.2382
	15					1.0000	0.9999	0.9979	0.9807	0.8818	0.5182
	16						1.0000	0.9998	0.9977	0.9775	0.8332
	17							1.0000	1.0000	1.0000	1.0000
18	0	0.1501	0.0180	0.0056	0.0016	0.0001	0.0000				
	1	0.4503	0.0991	0.0395	0.0142	0.0013	0.0001				
	2	0.7338	0.2713	0.1353	0.0600	0.0082	0.0007	0.0000			
	3	0.9018	0.5010	0.3057	0.1646	0.0328	0.0038	0.0002			
	4	0.9718	0.7164	0.5187	0.3327	0.0942	0.0154	0.0013	0.0000		
	5	0.9936	0.8671	0.7175	0.5344	0.2088	0.0481	0.0058	0.0003		
	6	0.9988	0.9487	0.8610	0.7217	0.3743	0.1189	0.0203	0.0014	0.0000	
	7	0.9998	0.9837	0.9431	0.8593	0.5634	0.2403	0.0576	0.0061	0.0002	
	8	1.0000	0.9957	0.9807	0.9404	0.7368	0.4073	0.1347	0.0210	0.0009	
	9		0.9991	0.9946	0.9790	0.8653	0.5927	0.2632	0.0596	0.0043	0.0000
	10		0.9998	0.9988	0.9939	0.9424	0.7597	0.4366	0.1407	0.0163	0.0002
	11		1.0000	0.9998	0.9986	0.9797	0.8811	0.6257	0.2783	0.0513	0.0012
	12			1.0000	0.9997	0.9942	0.9519	0.7912	0.4656	0.1329	0.0064
	13				1.0000	0.9987	0.9846	0.9058	0.6673	0.2838	0.0282
	14					0.9998	0.9962	0.9672	0.8354	0.4990	0.0982
	15					1.0000	0.9993	0.9918	0.9400	0.7287	0.2662
	16						0.9999	0.9987	0.9858	0.9009	0.5497
	17						1.0000	0.9999	0.9984	0.9820	0.8499
	18							1.0000	1.0000	1.0000	1.0000

<i>n</i>	<i>r</i>	<i>P</i>									
		0.10	0.20	0.25	0.30	0.40	0.50	0.60	0.70	0.80	0.90
19	0	0.1351	0.0144	0.0042	0.0011	0.0001					
	1	0.4203	0.0829	0.0310	0.0104	0.0008	0.0000				
	2	0.7054	0.2369	0.1113	0.0462	0.0055	0.0004	0.0000			
	3	0.8850	0.4551	0.2631	0.1332	0.0230	0.0022	0.0001			
	4	0.9648	0.6733	0.4654	0.2822	0.0696	0.0096	0.0006	0.0000		
	5	0.9914	0.8369	0.6678	0.4739	0.1629	0.0318	0.0031	0.0001		
	6	0.9983	0.9324	0.8251	0.6655	0.3081	0.0835	0.0116	0.0006		
	7	0.9997	0.9767	0.9225	0.8180	0.4878	0.1796	0.0352	0.0028	0.0000	
	8	1.0000	0.9933	0.9713	0.9161	0.6675	0.3238	0.0885	0.0105	0.0003	
	9		0.9984	0.9911	0.9674	0.8139	0.5000	0.1861	0.0326	0.0016	
	10		0.9997	0.9977	0.9895	0.9115	0.6762	0.3325	0.0839	0.0067	0.0000
	11		1.0000	0.9995	0.9972	0.9648	0.8204	0.5122	0.1820	0.0233	0.0003
	12			0.9999	0.9994	0.9884	0.9165	0.6919	0.3345	0.0676	0.0017
	13			1.0000	0.9999	0.9969	0.9682	0.8371	0.5261	0.1631	0.0086
	14				1.0000	0.9994	0.9904	0.9304	0.7178	0.3267	0.0352
	15					0.9999	0.9978	0.9770	0.8668	0.5449	0.1150
	16					1.0000	0.9996	0.9945	0.9538	0.7631	0.2946
	17						1.0000	0.9992	0.9896	0.9171	0.5797
	18							0.9999	0.9989	0.9856	0.8649
	19							1.0000	1.0000	1.0000	1.0000
20	0	0.1216	0.0115	0.0032	0.0008	0.0000					
	1	0.3917	0.0692	0.0243	0.0076	0.0005	0.0000				
	2	0.6769	0.2061	0.0913	0.0355	0.0036	0.0002				
	3	0.8670	0.4114	0.2252	0.1071	0.0160	0.0013	0.0000			
	4	0.9568	0.6296	0.4148	0.2375	0.0510	0.0059	0.0003			
	5	0.9887	0.8042	0.6172	0.4164	0.1256	0.0207	0.0016	0.0000		
	6	0.9976	0.9133	0.7858	0.6080	0.2500	0.0577	0.0065	0.0003		
	7	0.9996	0.9679	0.8982	0.7723	0.4159	0.1316	0.0210	0.0013	0.0000	
	8	0.9999	0.9900	0.9591	0.8867	0.5956	0.2517	0.0565	0.0051	0.0001	
	9	1.0000	0.9974	0.9861	0.9520	0.7553	0.4119	0.1275	0.0171	0.0006	
	10		0.9994	0.9961	0.9829	0.8725	0.5881	0.2447	0.0480	0.0026	0.0000
	11		0.9999	0.9991	0.9949	0.9435	0.7483	0.4044	0.1133	0.0100	0.0001
	12		1.0000	0.9998	0.9987	0.9790	0.8684	0.5841	0.2277	0.0321	0.0004
	13			1.0000	0.9997	0.9935	0.9423	0.7500	0.3920	0.0867	0.0024
	14				1.0000	0.9984	0.9793	0.8744	0.5836	0.1958	0.0113
	15					0.9997	0.9941	0.9490	0.7625	0.3704	0.0432
	16					1.0000	0.9987	0.9840	0.8929	0.5886	0.1330
	17						0.9998	0.9964	0.9645	0.7939	0.3231
	18						1.0000	0.9995	0.9924	0.9308	0.6083
	19							1.0000	0.9992	0.9885	0.8784
	20								1.0000	1.0000	1.0000

# The mean and variance of the Hypergeometric Distribution [1]

The mean and variance of the hypergeometric distribution  $h(x; N, n, k)$  are:

$$\text{Mean} = \frac{nk}{N}$$

$$\text{Variance} = \left( \frac{N-n}{N-1} \right) * \frac{nk}{N} * \left( \frac{N-k}{N} \right)$$

# The mean and variance of the Hypergeometric Distribution [2]

**Example:** Calculate the mean and variance of a hypergeometric random variable for parameters  $N = 700$ ,  $k = 35$ , and  $n = 20$ .

## Solution:

$$\text{Mean} = \frac{nk}{N}$$

$$\begin{aligned}\text{Mean} &= \frac{(20)(35)}{700} \\ &= 1\end{aligned}$$

$$\text{Variance} = \left( \frac{N-n}{N-1} \right) * \frac{nk}{N} * \left( \frac{N-k}{N} \right)$$

$$\begin{aligned}\text{Variance} &= \left( \frac{700-20}{700-1} \right) * \frac{(20)(35)}{700} * \left( \frac{700-35}{700} \right) \\ &= 0.9242\end{aligned}$$



# Relationship to the Binomial Distribution

## [1]

- ❑ There is an interesting relationship between the: **hypergeometric** and the **binomial distribution**. As one might expect, if **n is small compared to N**, the nature of the **N** items changes **very little** in each draw.
- ❑ So a **binomial distribution** can be used to approximate the **hypergeometric distribution** when **n is small, compared to N**.
- ❑ In fact, as a **rule of thumb** the **approximation** is good when  $\frac{n}{N} \leq 0.05$  or 5 %.

# Relationship to the Binomial Distribution [2]

- As a result, the **binomial distribution** may be viewed as a **large population** edition of the **hypergeometric distributions**

The mean and variance then come from the formulas

$$\text{Mean} = \mathbf{np} = \frac{nk}{N}$$

$$\text{Variance} = \mathbf{npq} = \frac{nk}{N} * \left(\frac{N-k}{N}\right)$$

$\frac{N-n}{N-1}$  is **negligible** when **n is small** relative to **N**

# Relationship to the Binomial Distribution [3]

**Example:** A manufacturer of automobile tires reports that among a shipment of **5000** sent to a local distributor, **1000** are slightly blemished. If one purchases **10** of these tires at random from the distributor, what, is the probability that exactly **3** are blemished?

## Solution:

Blemished	Non-blemished	Total
1000	4000	5000

Rule of thumb:  $\frac{n}{N} \leq 0.05 = \frac{10}{5000} = 0.002$  or 2% (true)

Here  $N = 5000$

$k = 1000$

$n = 10$

$p = k/N = 1000/5000 = 1/5 = 0.2$  (probability of blemished)

$X = 3$  (Let  $X$  denotes number of blemished tires)

$$h(3; 5000, 10, 1000) = b(3; 10, 0.2) = \sum_{x=0}^3 b(x; 10, 0.2) - \sum_{x=0}^2 b(x; 10, 0.2) = 0.8791 - 0.6778 = 0.2013.$$

# Hypergeometric Distribution

**Example:** Suppose that a shipment contains 5 defective items and 10 non defective items. If 7 items are selected at random without replacement, what is the probability that at least 3 defective items will be obtained?

## Solution:

Defective	Non defective	Total
5	10	15

Here  $N = 15$ ,  $n = 7$

$k = 5$  (defective items in the population)

Let  $X$  denotes number of defective items

$$P(X \geq 3) = 1 - P(X < 3) = 1 - \{P(X = 0) + P(X = 1) + P(X = 2)\}$$

$$h(x; N, n, k) = \frac{{}_k C_x {}_{N-k} C_{n-x}}{{}_N C_n},$$
$$\max\{0, n - (N - k)\} \leq x \leq \min\{n, k\}$$

$$\therefore \max\{0, n - (N - k)\} = \max\{0, 7 - (15 - 5)\} = 0$$

$$P(0) = \frac{{}^5C_0({}^{10}C_7)}{{}^{15}C_7} = 0.0186$$

$$P(1) = \frac{{}^5C_1({}^{10}C_6)}{{}^{15}C_7} = 0.1631$$

$$P(2) = \frac{{}^5C_2({}^{10}C_5)}{{}^{15}C_7} = 0.3916$$

$$P(X \geq 3) = 1 - P(X < 3) = 1 - (0.0186 + 0.1631 + 0.3916)$$

$$= 0.4267$$

**Example** A purchaser of electrical components buys them in **lots of size 10**. It is his policy to inspect **3 components** randomly from a lot and to **accept the lot** only if all **3 are nondefective**. If **30 percent** of the lots have **4 defective components** and **70 percent** have only **1**, what proportion of lots does the **purchaser reject**?



**30% of the lot:** Let  $x$  denotes the defective items from the **30% of the lot**.  $N = 10$ ,  $n = 3$ ,  $x = 0$ ,  **$k = 4$  (# of defectives items in 30% of the lot)**

$$h(x; N, n, k) = \frac{{}_k C_x ({}_{N-k} C_{n-x})}{{}_N C_n}$$

$$h(0; 10, 3, 4) = \frac{{}_4 C_0 ({}_6 C_3)}{{}_{10} C_3}$$

**70% of the lot :** Let  $y$  denotes the defective items from the remaining **70% of lot**.  $N = 10$ ,  $n = 3$ ,  $y = 0$ ,  **$k = 1$  (# of defectives items in 70% of the lot)**

$$h(x; N, n, k) = \frac{{}_k C_x ({}_{N-k} C_{n-x})}{{}_N C_n}$$

$$h(0; 10, 3, 1) = \frac{{}_1 C_0 ({}_9 C_3)}{{}_{10} C_3}$$

Let **A** denote the event that the purchaser **accepts a lot**.

$$\therefore P(A) = P(A \mid \text{lot has 4 defectives})\left(\frac{3}{10}\right) + P(A \mid \text{lot has 1 defective})\left(\frac{7}{10}\right)$$

$$P(A) = \frac{{}_4C_0 {}_6C_3}{{}_{10}C_3} \left(\frac{3}{10}\right) + \frac{{}_1C_0 {}_9C_3}{{}_{10}C_3} \left(\frac{7}{10}\right)$$

$$= \frac{54}{100} \text{ or } 54 \%$$

$$P(A^c) = 1 - P(A) = 1 - 0.54 = 0.46 \text{ or } 46\%$$

Hence, 46 percent of the lots are rejected.

# Mean and Variance for discrete probability distribution [1]

□ Expected value or mathematical expectation or **expectation** for discrete probability distribution is denoted by  $E(X)$  and is defined as

$$\begin{aligned} E(X) &= x_1 P(X_1 = x_1) + x_2 P(X_2 = x_2) + \dots + x_n P(X_n = x_n) \\ &= \sum_{i=1}^n x_i P(X_i = x_i) \end{aligned}$$

Here  $E(X)$  is mean or expected value of  $X$ .

# Mean and Variance for discrete probability distribution [2]

$$\begin{aligned} E(X^2) &= x_1^2 P(X_1 = x_1) + x_2^2 P(X_2 = x_2) + \dots + x_n^2 P(X_n = x_n) \\ &= \sum_{i=1}^n x_i^2 P(X_i = x_i) \end{aligned}$$

If  $X$  is a discrete random variable taking the values  $x_1, x_2, \dots, x_n$ , and having probability function  $P(x)$ , then the **variance** is given by

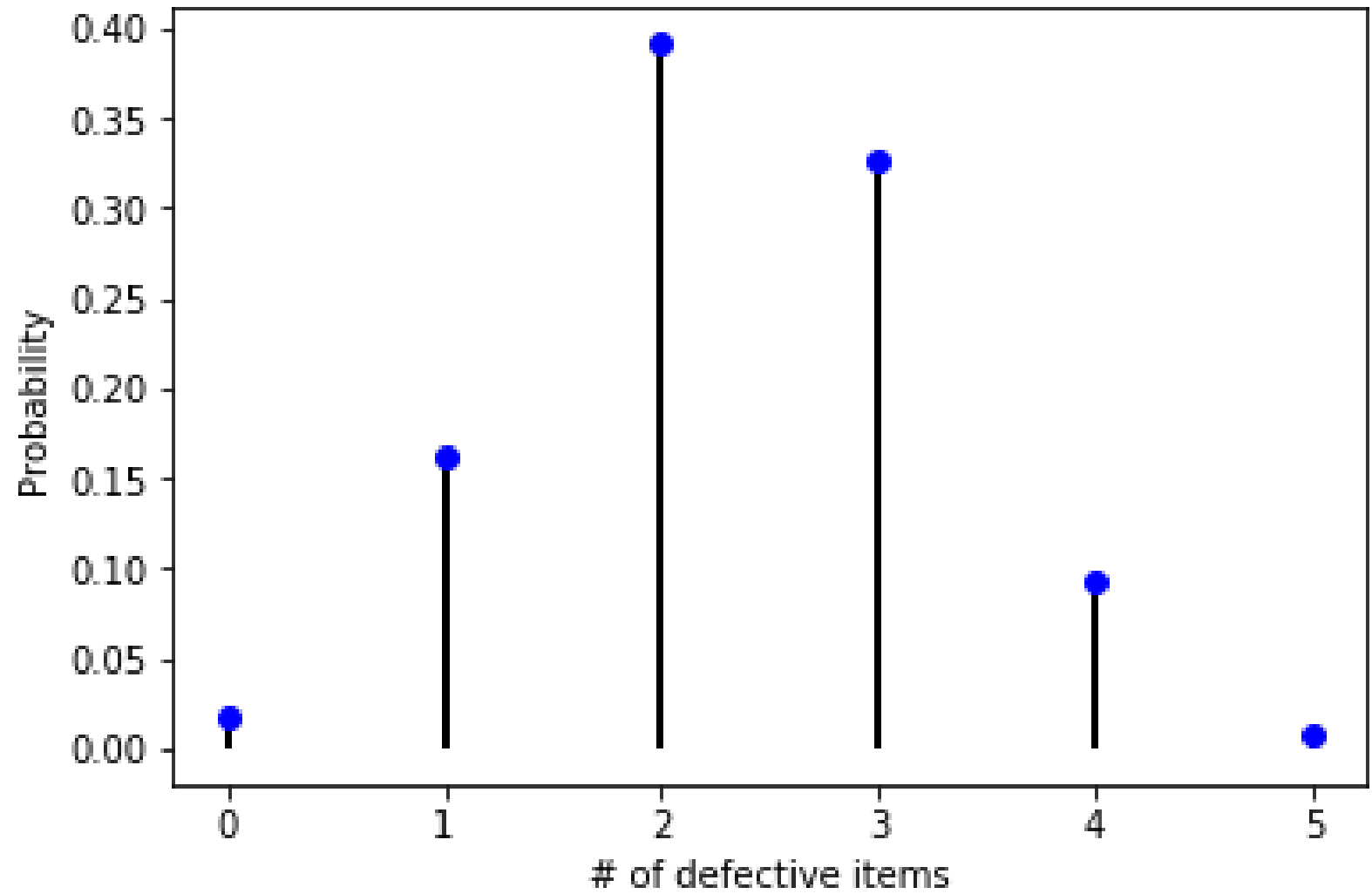
$$\text{Var}(X) = E(X^2) - \{E(X)\}^2$$

**Example** Suppose that a shipment contains **5 defective items** and 10 non defective items. If 7 items are selected at random without replacement, what is the **probability distribution of defective items**? Implement it in Python.

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import hypergeom
```

```
[N, n, k] = [15, 7, 5]
lLimit = max(0, n - (N - k))
uLimit = min(n, k)
rv = hypergeom(N, n, k)
x = np.arange(lLimit, uLimit + 1)
probability = rv.pmf(x)
```

```
fig = plt.figure()
ax = fig.add_subplot(111)
ax.plot(x, probability, 'bo')
ax.vlines(x, 0, probability, lw = 2)
ax.set_xlabel('# of defective items')
ax.set_ylabel('Probability')
plt.show()
```





**Example:** An industrial psychologist administered a personality inventory test for passive-aggressive traits to 150 employees. Each individual was given a score from 1 to 5, where 1 is extremely passive and 5 is extremely aggressive. A score of 3 indicated neither trait. The results are shown at the left. Construct a probability distribution for the random variable  $x$ . Find mean and variance of  $x$ .

### Frequency Distribution

Score, $x$	frequency, $f$
1	24
2	33
3	42
4	30
5	21

# Probability Distribution

$x$	$P(x)$	$xP(x)$	$x^2P(x)$
1	$\frac{24}{150} = 0.1600$	0.1600	0.1600
2	$\frac{33}{150} = 0.2200$	0.4400	0.8800
3	$\frac{42}{150} = 0.2800$	0.8400	2.5200
4	$\frac{30}{150} = 0.2000$	0.8000	3.2000
5	$\frac{21}{150} = 0.1400$	0.7000	3.5000
	$\sum \sum_{i=1}^5 \sum P_i = 1.0000$	$\sum xP(x) = 2.9400$	$\sum x^2P(x) = 10.2600$

**Note:**

**1.  $0 \leq P(x) \leq 1$**

**2.  $\sum P(x) = 1$**

$$\mu = E(x) = \sum xP(x) = 2.9400$$

$$E(x^2) = \sum x^2 P(x) = 10.2600$$

$$\sigma^2 = E(x^2) - [E(x)]^2 = 10.2600 - (2.9400)^2$$
$$\sigma^2 = 1.6164$$

# Multivariate Hypergeometric Distribution [1]

If  $N$  items can be partitioned into the  $k$  cells  $A_1, A_2, \dots, A_k$  with  $a_1, a_2, \dots, a_k$  elements, respectively, then the probability distribution of the random variables  $X_1, X_2, \dots, X_k$ , representing the number of elements selected from  $A_1, A_2, \dots, A_k$  in a random sample of size  $n$ , is

$$f(x_1, x_2, \dots, x_k; a_1, a_2, \dots, a_k, N, n) = \frac{\{({}^{a_1}C_{x_1}) ({}^{a_2}C_{x_2}) \dots ({}^{a_n}C_{x_n})\}}{{}_N C_n}$$

# Multivariate Hypergeometric Distribution [2]

**Example:** A group of **10** individuals is used for a biological case study. The group contains **3** people with blood type **O**, **4** with blood type **A**, and **3** with blood type **B**. What is the probability that a random sample of **5** will contain **1** person with blood type **O**, **2** people with blood type **A**, and **2** people with blood type **B**?

# Multivariate Hypergeometric Distribution [3]

**Solution :**  $a_1$  (type **O**) = 3,  $a_2$  = 4 (type **A**),  $a_3$  = 3 (type **B**)

$$x_1 = 1, x_2 = 2, x_3 = 2,$$

$$N = 10$$

$$n = 5$$

$$f(x_1, x_2, \dots, x_k; a_1, a_2, \dots, a_k, N, n) = \frac{\{({}^{a_1}C_{x_1}) ({}^{a_2}C_{x_2}) \dots ({}^{a_n}C_{x_n})\}}{{}_N C_n}$$

$$f(1, 2, 2; 3, 4, 3, 10, 5) = \frac{\{({}^3C_1) ({}^4C_2) ({}^3C_2)\}}{{}_{10}C_5} = 3/14 \\ = 0.2143$$