

STAT 88: Lecture 8

Contents

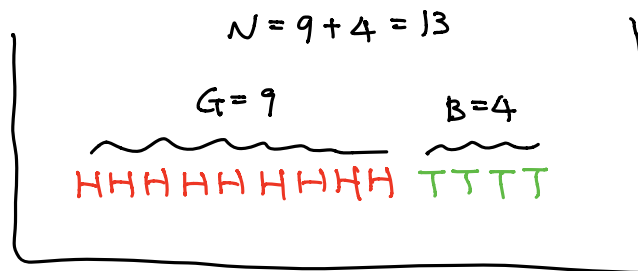
Section 3.5: Examples

Section 4.1: Cumulative Distribution Function (CDF)

Last time

Sec 3.4 The Hypergeometric distribution:

Picture:



↓ draw n w/o replacement

$g = 3, b = 2$ $n = g + b = 5$

HHH TT

Chance of g good in sample:

$$P(X = g) = \frac{\binom{G}{g} \binom{N-G}{n-g}}{\binom{N}{n}}$$

or in this case

$$\frac{\binom{9}{3} \binom{4}{2}}{\binom{13}{5}} = \binom{5}{3} \frac{9}{13} \frac{8}{12} \frac{7}{11} \frac{4}{10} \frac{3}{9}$$

Notice if draw w/ replacement
this is $\binom{5}{3} \left(\frac{9}{13}\right)^3 \left(\frac{4}{13}\right)^2$ "Binomial formula"

Warm up: Three cards are dealt from a standard 52 card deck. Find the chance:

- (a) The first card is red and the second two black.
- (b) Exactly one of the cards dealt is red.
- (c) At least one of the cards dealt is red.

3.5. Examples (continued)

Fisher Exact Test In a randomized controlled experiment with 100 participants, 60 participants are in the treatment group and 40 are in the control group. In the treatment group, 50 out of the 60 participants recover after the treatment. In the control group, 30 out of the 40 participants recover.

A total of 80 patients recovered out of 100.

Question. Suppose the treatment is not effective. What is the chance that 50 or more of the recovered patients are randomly assigned to the treatment group?

(if the answer is really small, then the treatment is probably effective.)

Start with: What is the chance that 50 of the recovered patients are randomly assigned to the treatment group?

4.1. Cumulative Distribution Function (CDF)

To specify a probability distribution, we have used a probability mass function (pmf):

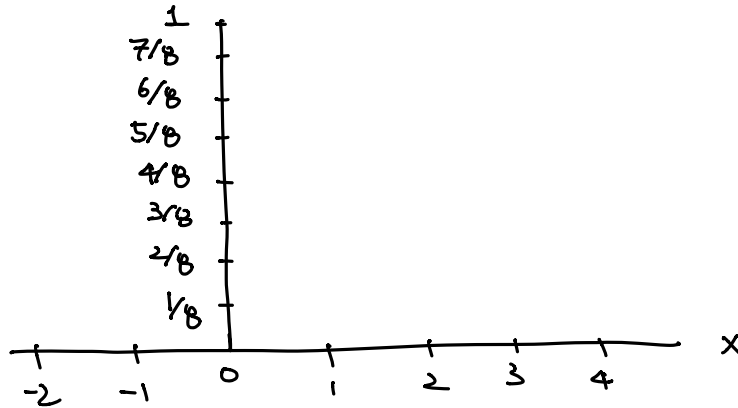
Example: $X \sim \text{Binomial}(3, 1/2)$.

You can also specify a probability distribution by giving the chance that the value of X is at most x , $F(x) = P(X \leq x)$. This is called the **cummulative distribution function (CDF)**.

Example:

x	0	1	2	3
$P(X=x)$	$1/8$	$3/8$	$3/8$	$1/8$
$F(x)$	$1/8$	$4/8$	$7/8$	1

Graph of the CDF We can define $F(x) = P(X \leq x)$ on the entire x-axis even though it only “jumps” at $x = 0, 1, 2, 3$.



Why does the CDF specify the distribution?

$$P(X = x) = P(X \leq x) - P(X \leq x - 1) = F(x) - F(x - 1).$$

So knowing $F(x)$ tells us $P(X = x)$.

Why is CDF useful? Solutions to many problems can be expressed in terms of CDF and Python has built-in CDF function.

Example: Fisher Exact test result.

$$1 - P(X < 50) = P(X \geq 50) = \sum_{g=50}^{60} \frac{\binom{80}{g} \binom{20}{60-g}}{\binom{100}{60}}.$$

" $1 - P(X \leq 49) = 1 - F(49)$

Computation You can use the stats module of SciPy to calculate CDF.

```
from scipy import stats
import numpy as np
```

```
1 - stats.hypergeom.cdf(49, 100, 80, 60)
```

```
0.22097998866696655
```

```
sum(stats.hypergeom.pmf(np.arange(50,61), 100, 80, 60))
```

```
0.22097998866696314
```

↖ old way

Example: In a population, 30% of the individuals are green and the rest are blue. Suppose you draw individuals with replacement until you draw a blue. Is the binomial formula applicable to find the chance that you draw 10 times?

Yes or no?