CS 237: Probability in Computing

Wayne Snyder Computer Science Department Boston University

Lecture 11:

- Random Variables (review)
- Probability Distributions of Random Variables
 - o Bernoulli

Negative Binomial (Pascal)

o **Binomial**

Hypergeometric

• Geometric

• [We will delay the Poisson]

Discrete RandomVariables: Probability Mass Function

The probability function of a **discrete** random variable X is a function

$$f_X$$
 = Probability Mass Function (PMF)

which assigns a probability to each real number in the range of X and follows the normal rules for a probability space:

$$f_{X} : R_{x} \to \mathcal{R}$$
$$\forall a \in R_{x} \quad f_{X}(a) \ge 0$$
$$\sum_{a \in R_{x}} f_{X}(a) = 1.0$$

Discrete RandomVariables: Probability Distributions

We will emphasize the distributions of random variables, using graphical representations to help our intuitions.

Example:

Y = "The number of tosses of a fair coin until a head appears"



Discrete RandomVariables

Notation:

$$P(X = k) =_{def} f_X(k)$$

$$P(X \neq k) =_{def} 1.0 - f_X(k)$$

$$P(X \leq k) =_{def} \sum_{a \leq k} f_X(a)$$

$$P(j \leq X \leq k) =_{def} \sum_{j \leq a \leq k} f_X(a)$$

$$P(Y = 4) = \frac{1}{16}$$

$$R_Y = \{1, 2, 3, ...\}$$

$$P(2 \le Y \le 4) = \frac{7}{16}$$



Probability Distribution for Y

Functions of Discrete RandomVariables

Why did I say you have to be careful? Two main reasons...

One, the function of a random variable may combine outcomes...

Example: Let Y' = X - 3 and let Y = |X - 3|

$$R_{Y'} = \{ -2, -1, 0, 1, 2, 3 \}$$
$$f_{Y'} = \{ \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6} \}$$

 $R_Y = \{0, 1, 2, 3\}$

 $f_{Y} = \{ \frac{1}{6}, \frac{2}{6}, \frac{2}{6}, \frac{1}{6} \}$



Functions of Discrete RandomVariables

$R_X = \{1, 2, 3, 4, 5, 6\}$ $f_X = \{\frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}\}$

Two, you have to be careful when a random variable is used more than once, since each occurrence refers to a potentially different random outcome!

Let Y = 2 * X (twice the dots showing on a thrown die)

 $R_Y = \{ 2, 4, 6, 8, 10, 12 \}$

 $\mathbf{f}_{\mathbf{V}} = \{ \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6} \}$



Let Y = X + X (sum of the dots showing on two thrown dice)





Standard Distributions

We will look at the following discrete distributions:

Discrete Distributions

- o Bernoulli
- o Binomial (Iterated version of Bernoulli)
- o Geometric
- Negative Binomial (Pascal) (Iterated version of Geometric)
- Hypergeometric (in your textbook, but we won't cover it)
- Poisson (We'll do this one separately in about a week).

These are summarized, with useful code for displaying the PMF and CDF in the notebook Distributions.ipynb on the class web site. Wikipedia has very good pages on all these distributions.

We will study these by considering the canonical problems which they describe....

Special Distributions

Any random variable X has a Probability Distribution which can be characterized by

- \circ R_X -- The range of the random variable
- \circ f_X -- The Probability Mass Function (PMF)
- \circ F_X -- The CDF
- \circ E(X) -- Expected value
- \circ Var(X), σ_X -- Variance, Standard Deviation

Next week...

In addition, we are interested in

- The canonical experiment which defines it
- \circ Formulae for calculating f_X and F_X , if such exist (hopefully efficient!)
- Algorithms for generating random variates from the distribution
- Any special properties of the distribution (e.g., the "memoryless property")
- Applications (random experiments which follow that distribution)

Bernoulli Distribution

Suppose you have a coin where the probability of a heads is p and we define the random variable

X = "the number of heads showing on one flipped coin"

Then we say that X is distributed according to the Bernoulli Distribution with parameter p, and write this as:

$$X \sim \text{Bernoulli}(p)$$

í

where



Jacob Bernoulli



	Jacob Bernoulli
Born	27 December 1654
	Basel, Switzerland
Died	16 August 1705 (aged 50)
	Basel, Switzerland

Among other accomplishments, Bernoulli discovered the number *e* (but Euler got the credit for "Euler's Number").

$$e = \lim_{n \to \infty} \left(1 + \frac{1}{n} \right)^n$$

Bernoulli Distribution

Each "poke" of such a random variable is called a Bernoulli Trial, and the outcomes are often labelled as

1 =Success 0 =Failure

Bernoulli Trials, and Bernoulli random variables, describe simple random experiments where there are two possible outcomes, and the probability of the outcomes is fixed by p and (1-p); the notion of "success" and "failure" is just a convenience and does not always correspond to the desirability of the outcome:

- Will I pass this course or not?
- Am I pregnant or not?
- Do I have cancer or not?



 $X \sim \text{Bernoulli}(p)$

 $R_X = \{0, 1\}$ $f_X = \{1 - p, p\}$



The Binomial Distribution occurs when you count the number of successes in N independent and identically distributed Bernoulli Trials (i.e., p is the same each time).

Formally, if Y ~ Bernoulli(p), and N times

$$X =$$
 "The number of successes in N trials of Y" = $Y + Y + ... + Y$

then we say that X is distributed according to the Binomial Distribution with parameters N and p, and write this as:

$$X \sim B(N,p)$$

Where

$$R_X = \{0, \dots, N\}$$
$$f_X(k) = \binom{N}{k} p^k (1-p)^{N-k}$$

Note: k successes and N-k failures SSS...FFFF... has probability $p^k (1-p)^{N-k}$ and there are such sequences.





For the visual display of this distribution, we will look briefly at the Distributions notebook to get a sense for this.....

The motivation for this distribution comes from the fact that many complex phenomena are composed of the additive effect of many small binary choices or events (Bernoulli Trials!); a vivid illustration of this can be seen in the Galton Board or Quicunx:

https://www.mathsisfun.com/data/quincunx.html

https://www.youtube.com/watch?v=J7AGOptcR1E

Example:

Suppose you have not studied for a True/False test and you randomly guess at each problem, and your probability of getting any particular problem correct is 50%; if there are 30 questions on the test and passing is a 60, what is your probability of passing?

$$P(X \ge 18) = \sum_{18 \le k \le 30} {30 \choose k} 0.5^k \cdot (1 - 0.5)^{30 - k}$$

In [7]: def f(N,p,k):
 return comb(N,k) * (p ** k) * ((1-p) ** (N-k))
sum([f(30,0.5,x) for x in range(18,31)])

```
Out[7]: 0.1807973040267825
```



Phenomena explained by the binomial are widespread throughout ordinary life, biology, engineering, and business:

- You go through 10 traffic lights; what is the probability that you stop at 4 of them?
- The probability of any individual in this class having a tattoo is 0.2; what is the probability that at least 40 people have a tattoo?
- Suppose 5% of tax returns are submitted with fraudulent data and the IRS examines 1% of returns; what is the probability that they will detect 3% of all fraudulent returns?
- There are about 700 gene variants which have been observed to have some influence on height; what is the probability that at least ³/₄'s of these genes will be dominant and have an influence on a person's height?

The binomial distribution is of widespread applicability, but it has a disadvantage: the only way to compute probabilities is to use the formula

$$P(X = k) = \binom{N}{k} p^k (1-p)^{N-k}$$

and this can involve some very large numbers.... for example:

Input:	
$\left(\begin{array}{c}10000\\5000\end{array}\right)$	
	$\binom{n}{m}$ is the binomial coefficient
Result:	Fewer digits More digits
$15917902635324389483375972736415211886530058374576145504283,\\ 191035177726371200957986632628539442222177433585982993226,\\ 205580463290870802073985087987219595848962041757866458580,\\ 184099587512068914331597813531740514534731996705213945025,\\ 384772773360083120537844882395127432175550288318092736464,\\ 428179545934936890023546288054736628292721322091972680306,\\ 215783976985524868345084786889499461126202336023529898945,\\ 892848842759111037432164623520292909554584530402349292778,\\ 714312397841036259290830007542173305536549242536830628153,\\ 072965334088925565069087515064761594462237620432685223006,\\ 267821125937595165711534284824533318106868409528400428469,\\ 950435925781799643074138942264944758662628186218375754128,\\ 036254688138854475912595618587146845438186146366235072846,\\ 821144165546574399328400579417002212869168618937974722788\\ \end{array}$	
Decimal approximation:	More digits
$\begin{array}{c} 1.591790263532438948337597273641521188653005837457614550428\ldots \times \\ 10^{3008} \end{array}$	
Number length:	
3009 decimal digits	

For comparison, there are about 10^{87} atoms in the universe....

The Geometric Distribution occurs when you count the number of independent and identically distributed Bernoulli trials until the first success.

Formally, if Y ~ Bernoulli(p), and

X = "The number of trials of Y until the first success"

then we say that X is distributed according to the Geometric Distribution with parameter p, and write this as: $X \sim G(p)$

where

$$R_X = \{ 1, 2, 3, \dots \}$$

 $f_X (k) = (1 - p)^{k-1} p$



For k, we have k-1 failures and 1 success (FFF... FS), which has probability (1-p)^{k-1} p.

Example

An absent-minded professor has 6 keys on his key ring and does not always remember which of his keys opens his office door. He chooses keys randomly and with replacement to try to open his door.

What is the probability that he opens it on the 3rd try?

Example

An absent-minded professor has 6 keys on his key ring and does not always remember which of his keys opens his office door. He chooses keys randomly and with replacement to try to open his door.

What is the probability that he opens it on the 3rd try?

Solution. This is G(1/6).

 $P(X=3) = (5/6)^2 (1/6) = 0.1157$

$R_X = \{$	1,	2,	3,	 <i>k</i> ,		}
$S = \{$	<i>S</i> ,	FS,	FFS,	 $FFF\ldots S$,		}
$f_X = \{$	<i>p</i> ,	(1-p)p,	$(1-p)^2p,$	 $(1-p)^{k-1}p,$	•••	}

But how do we know this is even a distribution? The only question is: Does f_X sum to 1.0?

Yes, no worries.... Suppose α is the sum of f_X . Then:

$$\alpha = p + (1 - p)p + (1 - p)^{2}p + (1 - p)^{3}p + \dots$$

= $p + (1 - p)(p + (1 - p)p + (1 - p)^{2}p + \dots)$
= $p + (1 - p)\alpha$
= $p + \alpha - p\alpha$

Subtracting α from both sides we have:

$$0 = p - p\alpha$$

$$\Leftrightarrow \qquad p\alpha = p$$

$$\Leftrightarrow \qquad \alpha = p/p = 1.0$$

Fortunately, the PMF is easy to compute, and there is are convenient formulae for inequalities:



$$R_X = \{ 1, 2, 3, \dots \}$$

 $f_X(k) = (1-p)^{k-1} p$

$$P(X > k) = (1-p)^{k} p + (1-p)^{k+1} p + (1-p)^{k+2} p + \dots$$

= $(1-p)^{k} (p + (1-p)p + (1-p)^{2}p + (1-p)^{3}p + \dots)$
= $(1-p)^{k}$

 $P_X(X \le k) \ = \ 1.0 - P_X(X > k) \ = \ 1.0 - (1 - p)^k$

Example

From an ordinary deck of 52 cards we draw cards at random and *with replacement*, and successively until an ace is drawn. What is the probability that at least 10 draws are needed?

Example

From an ordinary deck of 52 cards we draw cards at random and *with replacement*, and successively until an ace is drawn. What is the probability that at least 10 draws are needed?

Solution: The probability of an ace is 4/52 = 1/13. Thus:

 $P(X > 9) = (1 - 1/13)^9 = 0.4866$

Negative Binomial (Pascal) Distribution

The Negative Binomial is simply an "iterated" version of the Geometric.

Formally, if Y ~ Bernoulli(p) and

X = "The number of trials of Y until m successes occur"





then we say that X is distributed according to the Pascal Distribution with parameters m and p, and write:

$$X \sim Pascal(m, p)$$

where

$$R_X = \{m, m+1, m+2, ...\}$$

$$f_{X}(k) = {\binom{k-1}{m-1}} p^{m} (1-p)^{k-m}$$



Pascal Distribution

Example

Suppose you are throwing darts at a target for practice, and you decide that you will keep throwing until you hit the bull's eye 5 times. Suppose your probability of hitting the bullseye is 10%. What is the probability it takes exactly 10 throws?

Solution:

$$f_X(10) = \begin{pmatrix} 9 \\ 4 \end{pmatrix} 0.1^5 \cdot 0.9^5 = 7.44 \times 10^{-4}$$



C(9,4) * 0.1^(5) * 0.9 ^(5) ∫[™]₂₅ Extended Keyboard 1 Upload Assuming "C" is a math function | Use as a Input: $\binom{9}{4} \times 0.1^5 \times 0.9^5$ Result: 0.0007440174 18 # Just a more convenient syntax: 19 20 from scipy.special import comb 21 22 def C(N,K): 23 return comb(N,K,exact=True) 1 C(9,4)*0.1**5 * 0.9 ** 5

```
Out[6]: 0.000744017400000003
```

In [6]:

Pascal Distribution: Digression

CORRESPONDENCE

Open Access

Low dispersion in the infectiousness of COVID-19 cases implies difficulty in control



Daihai He^{1*}, Shi Zhao^{2,3}, Xiaoke Xu⁴, Qiangying Lin⁵, Zian Zhuang¹, Peihua Cao⁶, Maggie H. Wang^{2,3}, Yijun Lou¹, Li Xiao⁷, Ye Wu^{8,9*} and Lin Yang^{10*}

Abstract

The individual infectiousness of coronavirus disease 2019 (COVID-19), quantified by the number of secondary cases of a typical index case, is conventionally modelled by a negative-binomial (NB) distribution. Based on patient data of 9120 confirmed cases in China, we calculated the variation of the individual infectiousness, i.e., the dispersion parameter k of the NB distribution, at 0.70 (95% confidence interval: 0.59, 0.98). This suggests that the dispersion in the individual infectiousness is probably low, thus COVID-19 infection is relatively easy to sustain in the population and more challenging to control. Instead of focusing on the much fewer super spreading events, we also need to focus on almost every case to effectively reduce transmission.

Keywords: COVID-19, Basic reproductive number, Dispersion, Negative binomial, Mitigation



Article

On the Use of the Negative Binomial in Epidemiology

Prof. B. M. Bennett

First published: 1981 | https://doi.org/10.1002/bimj.4710230109 | Citations: 26

Cumulative Distribution Functions

One more topic before the lab tomorrow!

The Cumulative Distribution Function (CDF) for a random variable X shows what happens when we keep track of the sum of the probability distribution from left to right over its range:

$$F_X(k) = P(X \le k) = \sum_{a \le k} f_X(a)$$

Example: X = "The number of dots shown on a thrown die"

Probability Distribution Function P_X

Cumulative Distribution Function F_X



