

11: Joint (Multivariate) Distributions

Lisa Yan

April 29, 2020

Quick slide reference

3	Normal Approximation	11a_normal_approx
13	Discrete Joint RVs	11b_discrete_joint
26	Multinomial RV	11c_multinomial
34	Exercises	LIVE
43	Federalist Papers Demo	LIVE

Normal Approximation

Normal RVs

$$X \sim \mathcal{N}(\overset{\text{mean}}{\mu}, \overset{\text{variance}}{\sigma^2})$$

- Used to model many real-life situations because it maximizes entropy (i.e., randomness) for a given mean and variance
- Also useful for approximating the Binomial random variable!

Website testing

- 100 people are given a new website design.
- $X = \#$ people whose time on site increases
- The design actually has no effect, so $P(\text{time on site increases}) = 0.5$ independently.
- CEO will endorse the new design if $X \geq 65$.

What is $P(\text{CEO endorses change})$? Give a numerical approximation.

Approach 1: Binomial

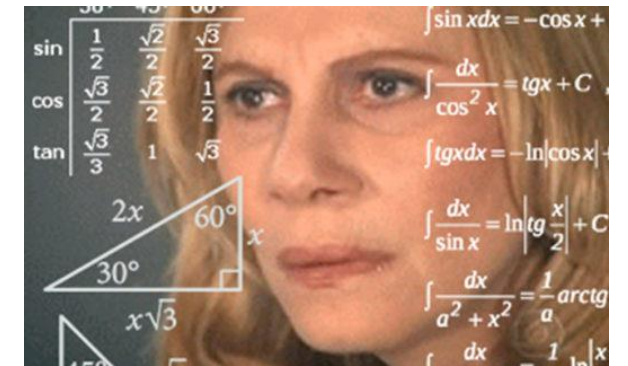
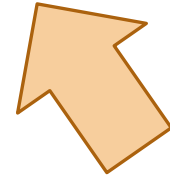
Define

$$X \sim \text{Bin}(n = 100, p = 0.5)$$

Want: $P(X \geq 65)$

Solve

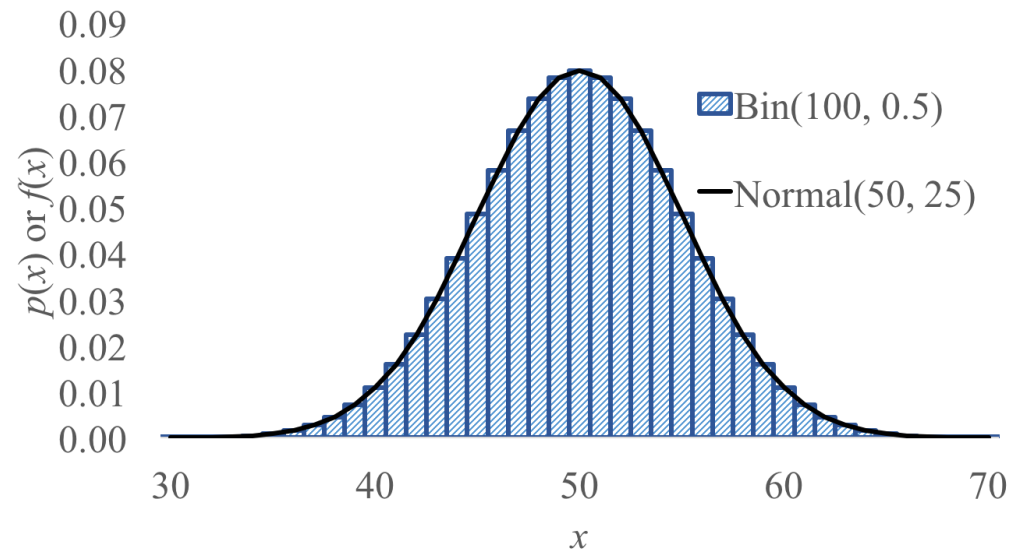
$$P(X \geq 65) = \sum_{i=65}^{100} \binom{100}{i} 0.5^i (1 - 0.5)^{100-i}$$



Don't worry, Normal approximates Binomial



Galton Board



(We'll explain *why*
in 2 weeks' time)

Website testing

- 100 people are given a new website design.
- $X = \#$ people whose time on site increases
- The design actually has no effect, so $P(\text{time on site increases}) = 0.5$ independently.
- CEO will endorse the new design if $X \geq 65$.

What is $P(\text{CEO endorses change})$? Give a numerical approximation.

Approach 1: Binomial

Define

$$X \sim \text{Bin}(n = 100, p = 0.5)$$

Want: $P(X \geq 65)$

Solve

$$P(X \geq 65) \approx \mathbf{0.0018}$$

Approach 2: approximate with Normal

Define

$$Y \sim \mathcal{N}(\mu, \sigma^2)$$

$$\mu = np = 50$$

$$\sigma^2 = np(1 - p) = 25$$

$$\sigma = \sqrt{25} = 5$$

Solve

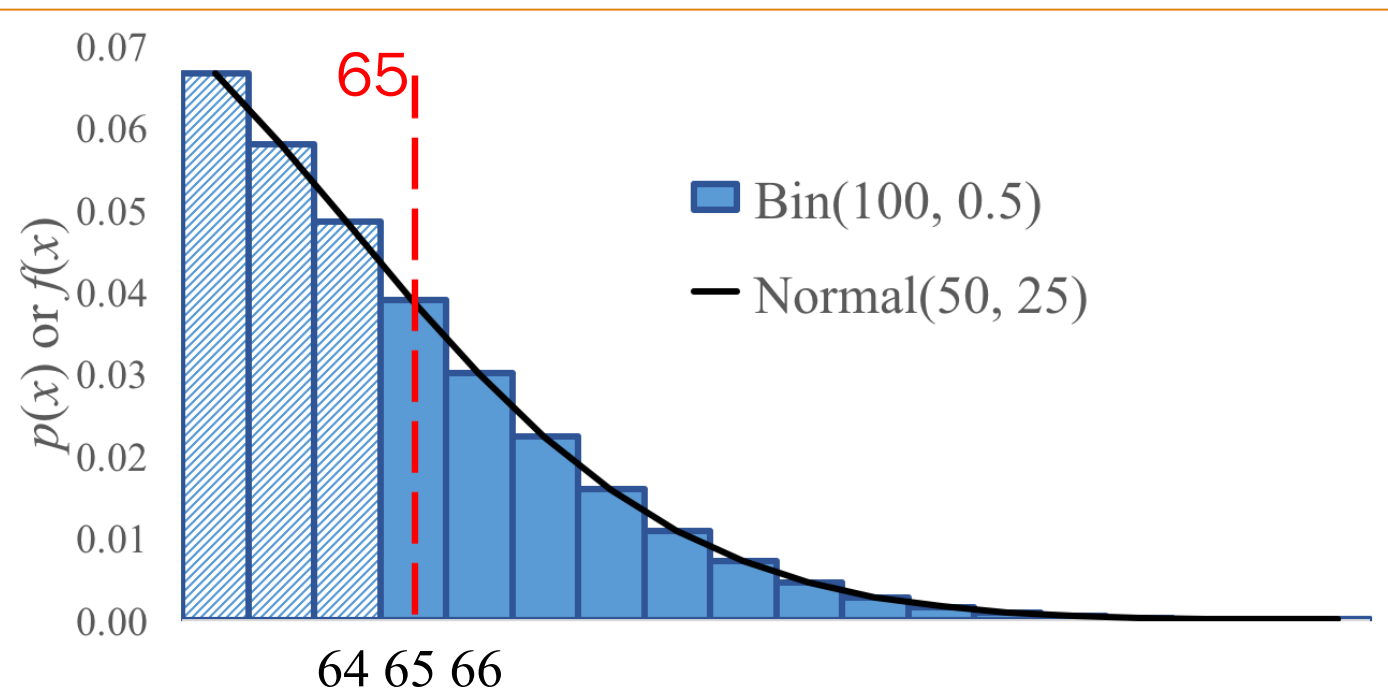
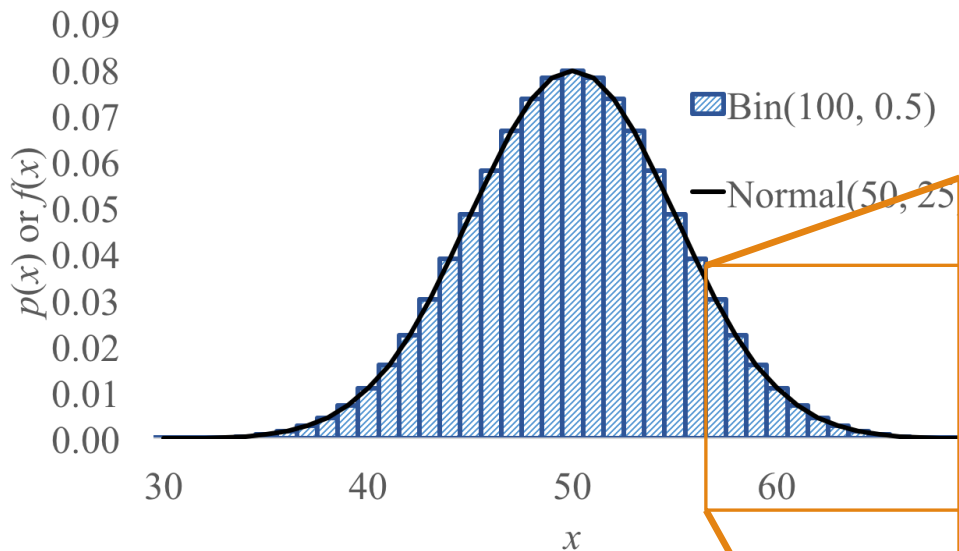
$$\begin{aligned} P(X \geq 65) &\approx P(Y \geq 65) = 1 - F_Y(65) \\ &= 1 - \Phi\left(\frac{65-50}{5}\right) = 1 - \Phi(3) \approx \mathbf{0.0013} \end{aligned}$$



(this approach is actually missing something)

Website testing (with continuity correction)

In our website testing, $Y \sim \mathcal{N}(50, 25)$ approximates $X \sim \text{Bin}(100, 0.5)$.



$$P(X \geq 65) \text{ Binomial}$$

$$\approx P(Y \geq 64.5) \text{ Normal}$$

$$\approx 0.0018 \quad \checkmark \text{ the better Approach 2}$$

You must perform a **continuity correction** when approximating a Binomial RV with a Normal RV.

Continuity correction

If $Y \sim \mathcal{N}(np, np(1 - p))$ approximates $X \sim \text{Bin}(n, p)$, how do we approximate the following probabilities?

Discrete (e.g., Binomial)
probability question



Continuous (Normal)
probability question

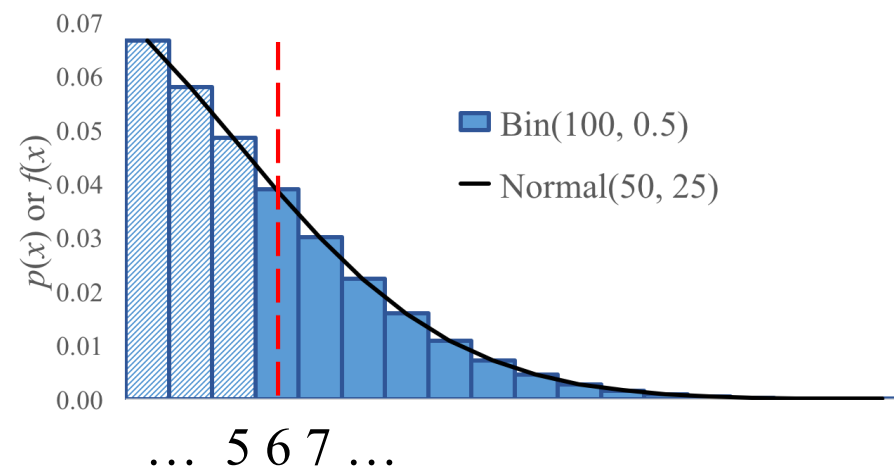
$$P(X = 6)$$

$$P(X \geq 6)$$

$$P(X > 6)$$

$$P(X < 6)$$

$$P(X \leq 6)$$



Continuity correction

If $Y \sim \mathcal{N}(np, np(1 - p))$ approximates $X \sim \text{Bin}(n, p)$, how do we approximate the following probabilities?

Discrete (e.g., Binomial)
probability question



Continuous (Normal)
probability question

$$P(X = 6)$$

$$P(X \geq 6)$$

$$P(X > 6)$$

$$P(X < 6)$$

$$P(X \leq 6)$$

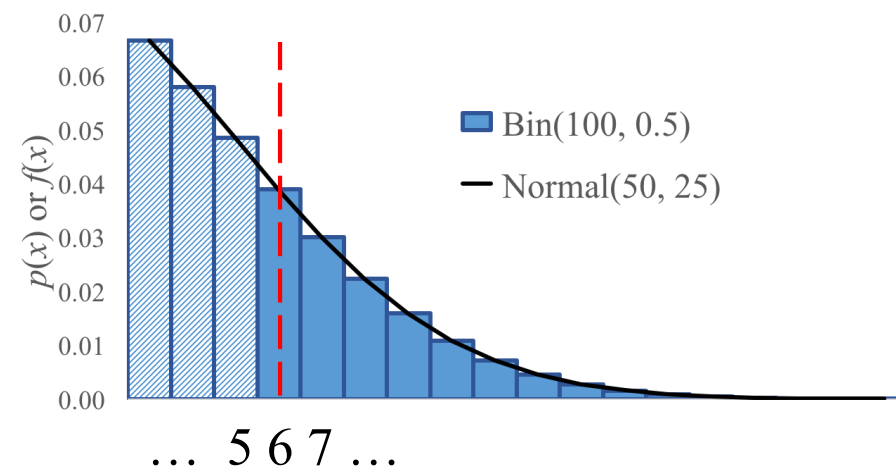
$$P(5.5 \leq Y \leq 6.5)$$

$$P(Y \geq 5.5)$$

$$P(Y \geq 6.5)$$

$$P(Y \leq 5.5)$$

$$P(Y \leq 6.5)$$



Who gets to approximate?

$$X \sim \text{Bin}(n, p)$$

$$E[X] = np$$

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



$$Y \sim \text{Poi}(\lambda)$$

$$\lambda = np$$

?

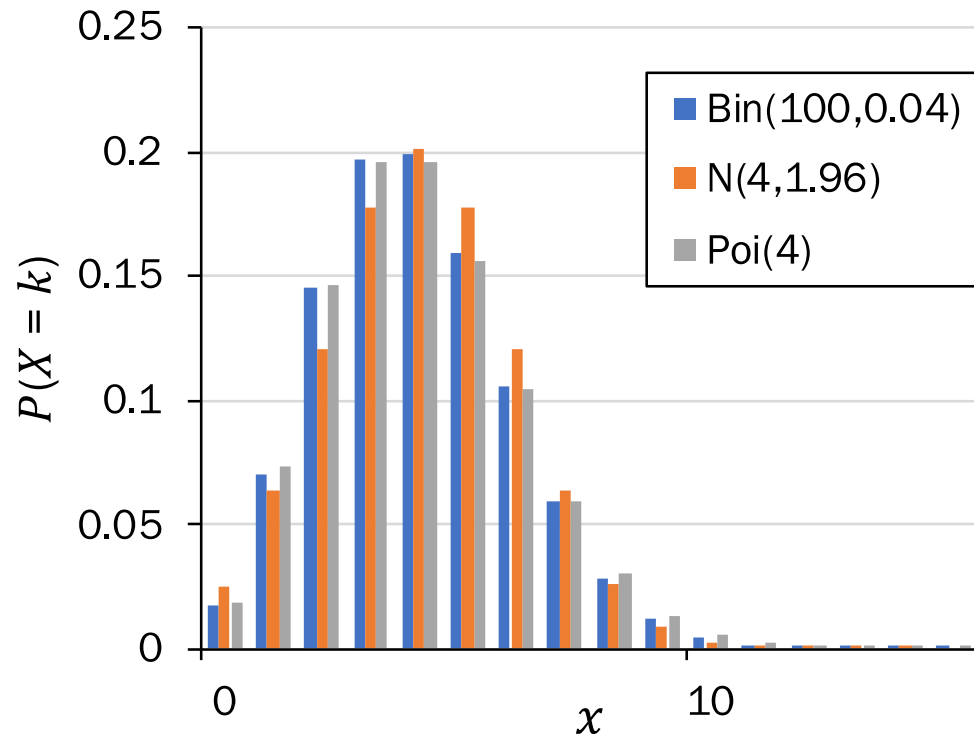


$$Y \sim \mathcal{N}(\mu, \sigma^2)$$

$$\mu = np$$

$$\sigma^2 = np(1 - p)$$

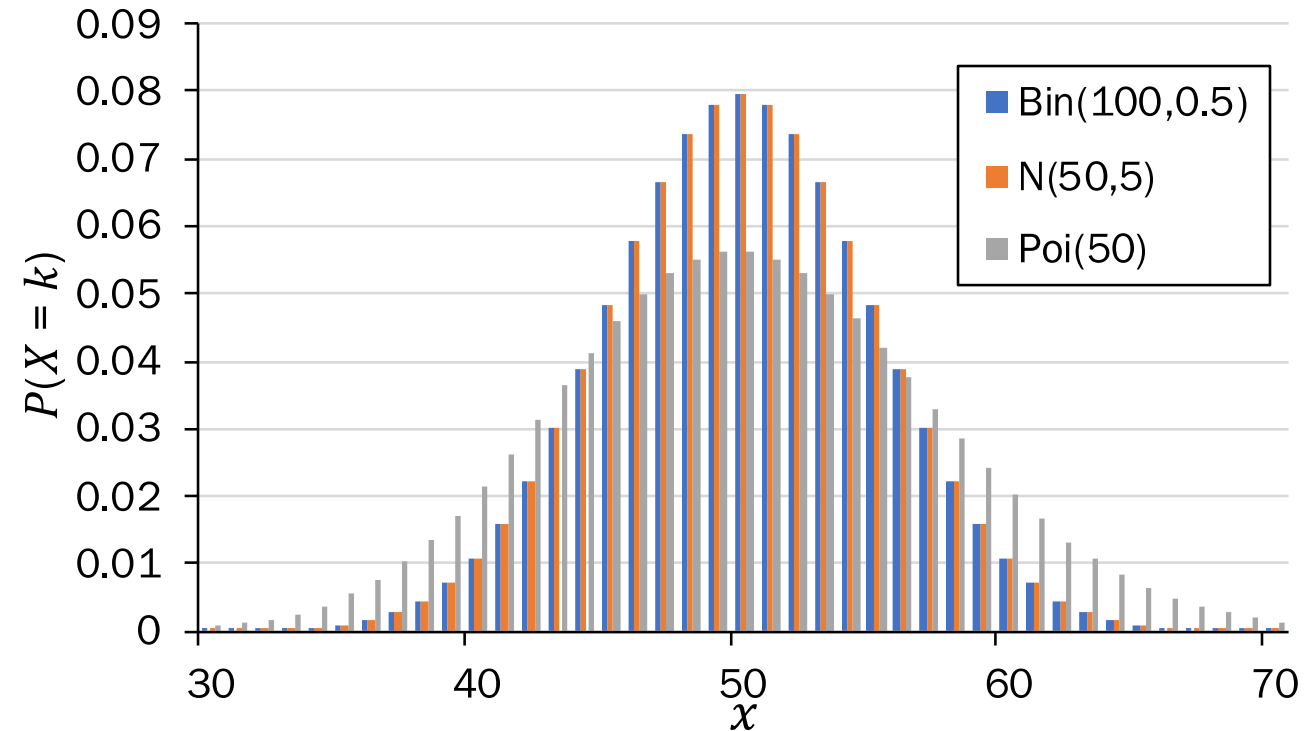
Who gets to approximate?



Poisson approximation

n large (> 20), p small (< 0.05)

slight dependence okay



Normal approximation

n large (> 20), p mid-ranged ($np(1 - p) > 10$)

independence

1. If there is a choice, use Normal to approximate.
2. When using Normal to approximate a discrete RV, use a continuity correction.

Discrete Joint RVs



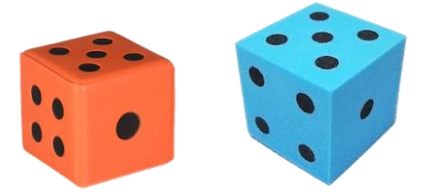
$$P(A_W > A_B)$$

This is a probability of an event involving *two* random variables!

What is the probability that the Warriors win?
How do you model zero-sum games?

Joint probability mass functions

Roll two 6-sided dice, yielding values X and Y .



X

random variable

$$P(X = 1)$$

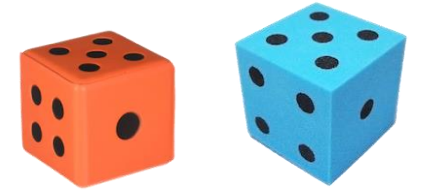
probability of
an event

$$P(X = k)$$

probability mass function

Joint probability mass functions

Roll two 6-sided dice, yielding values X and Y .

 X

random variable

$$P(X = 1)$$

probability of
an event

$$P(X = k)$$

probability mass function

 X, Y

random variables

$$P(X = 1 \cap Y = 6)$$

$$P(X = 1, Y = 6)$$

new notation: the comma

probability of the intersection
of two events

$$P(X = a, Y = b)$$

joint probability mass function

Discrete joint distributions

For two discrete joint random variables X and Y , the **joint probability mass function** is defined as:

$$p_{X,Y}(a, b) = P(X = a, Y = b)$$

The **marginal distributions** of the joint PMF are defined as:

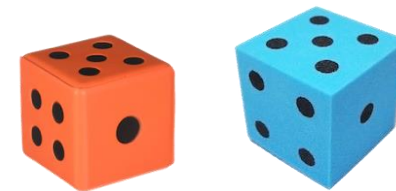
$$p_X(a) = P(X = a) = \sum_y p_{X,Y}(a, y)$$

$$p_Y(b) = P(Y = b) = \sum_x p_{X,Y}(x, b)$$

Use marginal distributions to get a 1-D RV from a joint PMF.

Two dice

Roll two 6-sided dice, yielding values X and Y .



1. What is the joint PMF of X and Y ?

$$p_{X,Y}(a, b) = 1/36 \quad (a, b) \in \{(1,1), \dots, (6,6)\}$$

		X					
		1	2	3	4	5	6
Y	1	1/36	1/36
	2
	3
	4
	5
	6	1/36	1/36

An orange arrow points from the text $P(X = 4, Y = 2)$ to the cell at $X=4, Y=2$ in the table.

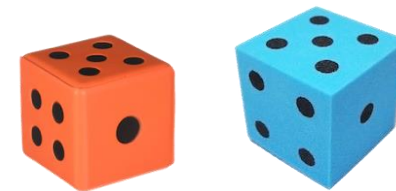
Probability table

- All possible outcomes for several discrete RVs
- Not parametric (e.g., parameter p in $\text{Ber}(p)$)

Two dice

Roll two 6-sided dice, yielding values X and Y .

1. What is the joint PMF of X and Y ?



$$p_{X,Y}(a, b) = 1/36 \quad (a, b) \in \{(1,1), \dots, (6,6)\}$$

2. What is the marginal PMF of X ?

$$p_X(a) = P(X = a) = \sum_y p_{X,Y}(a, y) = \sum_{y=1}^6 \frac{1}{36} = \frac{1}{6} \quad a \in \{1, \dots, 6\}$$

A computer (or three) in every house.

Consider households in Silicon Valley.

- A household has X Macs and Y PCs.
- Each house has a maximum of 3 computers (Macs + PCs) in the house.

1. What is $P(X = 1, Y = 0)$, the missing entry in the probability table?

		X (# Macs)			
		0	1	2	3
Y (# PCs)	0	.16	?	.07	.04
	1	.12	.14	.12	0
	2	.07	.12	0	0
	3	.04	0	0	0



A computer (or three) in every house.

Consider households in Silicon Valley.

- A household has X Macs and Y PCs.
- Each house has a maximum of 3 computers (Macs + PCs) in the house.

1. What is $P(X = 1, Y = 0)$, the missing entry in the probability table?

		X (# Macs)			
		0	1	2	3
Y (# PCs)	0	.16	.12	.07	.04
	1	.12	.14	.12	0
	2	.07	.12	0	0
	3	.04	0	0	0

A joint PMF must sum to 1:

$$\sum_x \sum_y p_{X,Y}(x, y) = 1$$

A computer (or three) in every house.

Consider households in Silicon Valley.

- A household has X Macs and Y PCs.
- Each house has a maximum of 3 computers (Macs + PCs) in the house.

2. How do you compute the marginal PMF of X ?

		X (# Macs)				
		0	1	2	3	
Y (# PCs)	0 A	.16	.12	.07	.04	.39
	1	.12	.14	.12	0	.38
	2	.07	.12	0	0	.19
	3	.04	0	0	0	.04
B		.39	.38	.19	.04	sum rows here



A computer (or three) in every house.

Consider households in Silicon Valley.

- A household has X Macs and Y PCs.
- Each house has a maximum of 3 computers (Macs + PCs) in the house.

2. How do you compute the marginal PMF of X ?

		X (# Macs)				
		0	1	2	3	
Y (# PCs)	0	.16	.12	.07	.04	.39
	1	.12	.14	.12	0	.38
	2	.07	.12	0	0	.19
	3	.04	0	0	0	.04
		.39	.38	.19	.04	sum cols here

A. $p_{X,Y}(x, 0) = P(X = x, Y = 0)$

B. Marginal PMF of X $p_X(x) = \sum_y p_{X,Y}(x, y)$

C. Marginal PMF of Y $p_Y(y) = \sum_x p_{X,Y}(x, y)$

To find a marginal distribution over one variable, sum over all other variables in the joint PMF.

A computer (or three) in every house.

Consider households in Silicon Valley.

- A household has X Macs and Y PCs.
- Each house has a maximum of 3 computers (Macs + PCs) in the house.

3. Let $C = X + Y$. What is $P(C = 3)$?

		X (# Macs)			
		0	1	2	3
Y (# PCs)	0	.16	.12	.07	.04
	1	.12	.14	.12	0
	2	.07	.12	0	0
	3	.04	0	0	0



A computer (or three) in every house.

Consider households in Silicon Valley.

- A household has X Macs and Y PCs.
- Each house has a maximum of 3 computers (Macs + PCs) in the house.

3. Let $C = X + Y$. What is $P(C = 3)$?

		X (# Macs)			
		0	1	2	3
Y (# PCs)	0	.16	.12	.07	.04
	1	.12	.14	.12	0
	2	.07	.12	0	0
	3	.04	0	0	0

$$P(C = 3) = P(X + Y = 3)$$

Law of Total Probability

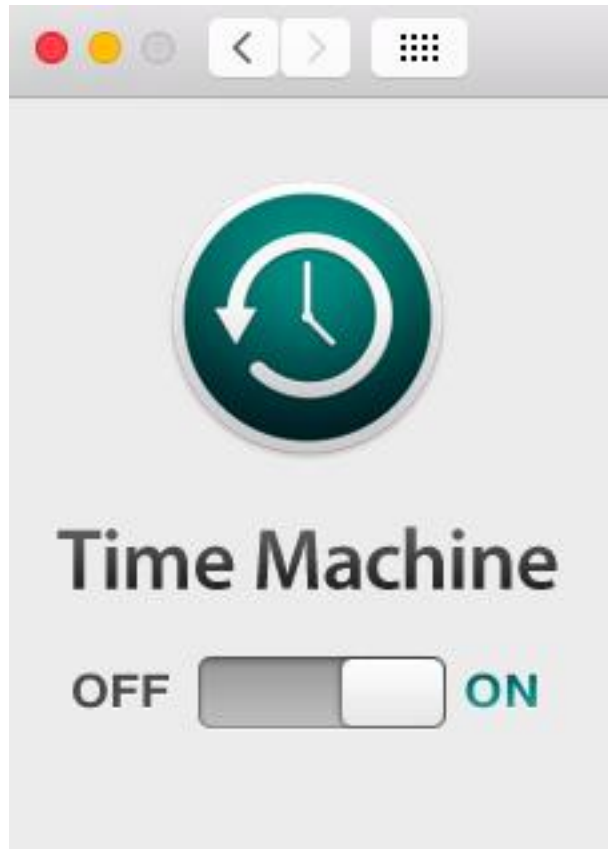
$$= \sum_x \sum_y P(X + Y = 3 | X = x, Y = y) P(X = x, Y = y)$$

$$= P(X = 0, Y = 3) + P(X = 1, Y = 2) \\ + P(X = 2, Y = 1) + P(X = 3, Y = 0)$$

We'll come back to sums of RVs next lecture!

Multinomial RV

Recall the good times



Permutations

$n!$

How many ways are there to order n objects?

Counting unordered objects

Binomial coefficient

How many ways are there to group n objects into **two** groups of size k and $n - k$, respectively?

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

Called the binomial coefficient because of something from Algebra

Multinomial coefficient

How many ways are there to group n objects into r groups of sizes n_1, n_2, \dots, n_r respectively?

$$\binom{n}{n_1, n_2, \dots, n_r} = \frac{n!}{n_1! n_2! \cdots n_r!}$$

Multinomials generalize Binomials for counting.

Probability

Binomial RV

What is the probability of getting k successes and $n - k$ failures in n trials?

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

Binomial # of ways of ordering the successes

Probability of each ordering of k successes is equal + mutually exclusive

Multinomial RV

What is the probability of getting c_1 of outcome 1, c_2 of outcome 2, ..., and c_m of outcome m in n trials?

Multinomial RVs also generalize Binomial RVs for probability!

Multinomial Random Variable

Consider an experiment of n independent trials:

- Each trial results in one of m outcomes. $P(\text{outcome } i) = p_i$, $\sum_{i=1}^m p_i = 1$
- Let $X_i = \#$ trials with outcome i

Joint PMF

$$P(X_1 = c_1, X_2 = c_2, \dots, X_m = c_m) = \binom{n}{c_1, c_2, \dots, c_m} p_1^{c_1} p_2^{c_2} \dots p_m^{c_m}$$

where $\sum_{i=1}^m c_i = n$ and $\sum_{i=1}^m p_i = 1$

Multinomial # of ways of ordering the outcomes

Probability of each ordering is equal + mutually exclusive

Hello dice rolls, my old friends

A 6-sided die is rolled 7 times.

What is the probability of getting:

- 1 one
- 0 threes
- 0 fives
- 1 two
- 2 fours
- 3 sixes



Hello dice rolls, my old friends

A 6-sided die is rolled 7 times.

What is the probability of getting:

- 1 one
- 1 two
- 0 threes
- 2 fours
- 0 fives
- 3 sixes

$$P(X_1 = 1, X_2 = 1, X_3 = 0, X_4 = 2, X_5 = 0, X_6 = 3)$$

$$= \binom{7}{1,1,0,2,0,3} \left(\frac{1}{6}\right)^1 \left(\frac{1}{6}\right)^1 \left(\frac{1}{6}\right)^0 \left(\frac{1}{6}\right)^2 \left(\frac{1}{6}\right)^0 \left(\frac{1}{6}\right)^3 = 420 \left(\frac{1}{6}\right)^7$$

Hello dice rolls, my old friends

A 6-sided die is rolled 7 times.

What is the probability of getting:

- 1 one • 0 threes • 0 fives
- 1 two • 2 fours • 3 sixes

of times
a six appears

$$P(X_1 = 1, X_2 = 1, X_3 = 0, X_4 = 2, X_5 = 0, X_6 = 3)$$

$$= \binom{7}{1,1,0,2,0,3} \left(\frac{1}{6}\right)^1 \left(\frac{1}{6}\right)^1 \left(\frac{1}{6}\right)^0 \left(\frac{1}{6}\right)^2 \left(\frac{1}{6}\right)^0 \left(\frac{1}{6}\right)^3 = 420 \left(\frac{1}{6}\right)^7$$

choose where
the sixes appear

probability
of rolling a six this many times

11: Joint (Multivariate) Distributions (live)

Slides by Lisa Yan
April 29, 2020

Normal RVs

$$X \sim \mathcal{N}(\overset{\text{mean}}{\mu}, \overset{\text{variance}}{\sigma^2})$$

- Used to model many real-life situations because it maximizes entropy (i.e., randomness) for a given mean and variance
- Also useful for approximating the Binomial random variable!

Who gets to approximate?

$$X \sim \text{Bin}(n, p)$$

$$E[X] = np$$

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



$$Y \sim \text{Poi}(\lambda)$$

$$\lambda = np$$

n large (> 20)

p small (< 0.05)

slight dependence okay

$$Y \sim \mathcal{N}(\mu, \sigma^2)$$

$$\mu = np$$

$$\sigma^2 = np(1 - p)$$

n large (> 20), p mid-ranged ($np(1 - p) > 10$)

independence

need continuity correction

- Computing probabilities on Binomial RVs is often computationally expensive.
- Two reasonable approximations, but when to use which?

Think

Check out the question on the next slide.
Post any clarifications here!

<https://us.edstem.org/courses/667/discussion/90049>



Stanford Admissions (a while back)

Stanford accepts 2480 students.

- Each accepted student has 68% chance of attending (independent trials)
- Let $X = \#$ of students who will attend

What is $P(X > 1745)$? *Give a numerical approximation.*

- Strategy:
- A. Just Binomial
 - B. Poisson
 - C. Normal
 - D. None/other



Stanford Admissions (a while back)

Stanford accepts 2480 students.

- Each accepted student has 68% chance of attending (independent trials)
- Let $X = \#$ of students who will attend

What is $P(X > 1745)$? Give a numerical approximation.

Strategy: A. Just Binomial not an approximation (also computationally expensive)

B. Poisson $p = 0.68$, not small enough

C. Normal Variance $np(1 - p) = 540 > 10$

D. None/other

Define an approximation

Let $Y \sim \mathcal{N}(E[X], \text{Var}(X))$

$$E[X] = np = 1686$$

$$\text{Var}(X) = np(1 - p) \approx 540 \rightarrow \sigma = 23.3$$

$$P(X > 1745) \approx P(Y \geq 1745.5) \quad \triangle! \text{ Continuity correction}$$

Lisa Yan, CS109, 2020

Solve

$$P(Y \geq 1745.5) = 1 - F(1745.5)$$

$$= 1 - \Phi\left(\frac{1745.5 - 1686}{23.3}\right)$$

$$= 1 - \Phi(2.54) \approx 0.0055$$

Changes in Stanford Admissions

Stanford accepts 2480 students.

- Each accepted student has 68% chance of attending (independent trials)
- Let $X = \#$ of students who will attend

Yield rate 20
years ago

What is $P(X > 1745)$? Give a numerical approximation.



The Stanford Daily

NEWS · SPORTS · OPINIONS · ARTS & LIFE · THE GRIND · MULTIMEDIA · FEATURES · ARCHIVES

Class of 2018 admit rates lowest in University history

March 28, 2014 16 Comments [Tweet](#) [Like 901](#)

Alex Zivkovic
Desk Editor

Stanford admitted 2,138 students to the Class of 2018 in this year's admissions cycle, producing – at 5.07 percent – the lowest admit rate in University history.

The [University](#) received a total of 42,167 applications this year, a record total and a 8.6 percent increase over [last year's figure of 38,828](#). Stanford [accepted 748 students](#)



Overview for the Class of 2022

- Total Applicants: 47,451 Admit rate: 4.3%
- Total Admits: 2,071 Yield rate: 81.9%
- Total Enrolled: 1,706

People love coming to Stanford!

Consider an experiment of n independent trials:

- Each trial results in one of m outcomes. $P(\text{outcome } i) = p_i$, $\sum_{i=1}^m p_i = 1$
- Let $X_i = \#$ trials with outcome i

Joint PMF

$$P(X_1 = c_1, X_2 = c_2, \dots, X_m = c_m) = \binom{n}{c_1, c_2, \dots, c_m} p_1^{c_1} p_2^{c_2} \cdots p_m^{c_m}$$

where $\sum_{i=1}^m c_i = n$ and $\sum_{i=1}^m p_i = 1$

Example:

- Rolling 2 twos, 3 threes, and 5 fives on 10 rolls of a fair-sided die
- Generating a random 5-word phrase with 1 “the”, 2 “bacon”, 1 “put”, 1 “on”

Hello dice rolls, my old friends

A 6-sided die is rolled 7 times.

What is the probability of getting:

- 1 one
- 1 two
- 0 threes
- 2 fours
- 0 fives
- 3 sixes

of times
a six appears

$$P(X_1 = 1, X_2 = 1, X_3 = 0, X_4 = 2, X_5 = 0, X_6 = 3)$$

$$= \binom{7}{1,1,0,2,0,3} \left(\frac{1}{6}\right)^1 \left(\frac{1}{6}\right)^1 \left(\frac{1}{6}\right)^0 \left(\frac{1}{6}\right)^2 \left(\frac{1}{6}\right)^0 \left(\frac{1}{6}\right)^3 = 420 \left(\frac{1}{6}\right)^7$$

choose where
the sixes appear

probability
of rolling a six this many times

Parameters of a Multinomial RV?

$X \sim \text{Bin}(n, p)$ has parameters $n, p \dots$

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

p : probability of success outcome on a single trial

A Multinomial RV has parameters n, p_1, p_2, \dots, p_m (Note $p_m = 1 - \sum_{i=1}^{m-1} p_i$)

$$P(X_1 = c_1, X_2 = c_2, \dots, X_m = c_m) = \binom{n}{c_1, c_2, \dots, c_m} p_1^{c_1} p_2^{c_2} \cdots p_m^{c_m}$$

p_i : probability of outcome i on a single trial

Where do we get p_i from?

Interlude for fun/announcements

Announcements

More OH!

Interesting probability news

Estimating Coronavirus Prevalence by Cross-Checking Countries

<https://medium.com/@jsteinhardt/estimating-coronavirus-prevalence-by-cross-checking-countries-c7e4211f0e18>

We'll make the modeling assumption that N_{ij} is a **Poisson distribution with rate parameter** $A_{ij} * \lambda_i * \alpha_j$. What this means is that the **expected number of cases** should be equal to the total amount of travel, times some source-dependent multiplier α_j ..., times some country-dependent multiplier λ_i (the infection prevalence in country i).”



Ethics in Probability: Biased Data + Bayes

Amazon scraps secret AI recruiting tool that showed bias against women

“In effect, Amazon’s system *taught itself that male candidates were preferable*. It penalized resumes that included the word ‘women’s,’ as in ‘women’s chess club captain.’”

Basic Bayes Algorithm: Pick highest $P(H | \text{resume})$

Let resume_F be a resume associated with a female applicant.

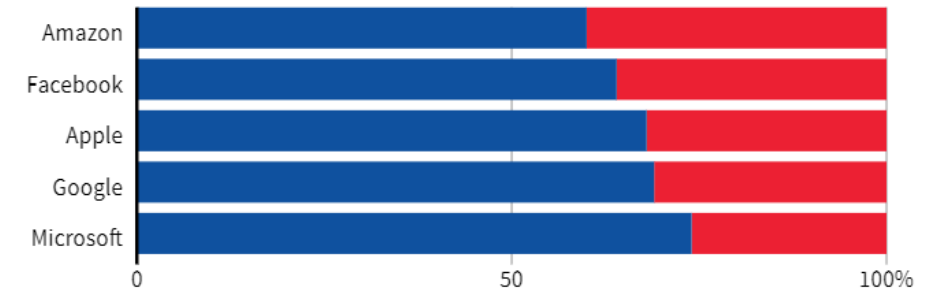
$$P(H | \text{resume}_F) = P(\text{resume}_F | H) * P(H) / P(\text{resume}_F)$$

Because of biased historical data, $P(\text{resume}_F | H)$ is small ☹️ ➔

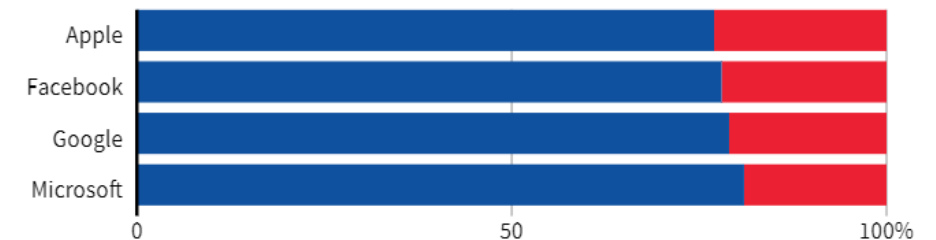
Therefore $P(H | \text{resume}_F)$ may be higher than $P(H | \text{resume}_M)$ *simply because of biased historical data, rather than comparative candidate skillsets.*

GLOBAL HEADCOUNT

■ Male ■ Female



EMPLOYEES IN TECHNICAL ROLES



<https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scraps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G>

Ethics in Probability: Biased Data + Bayes

What if we ignore gender traits?

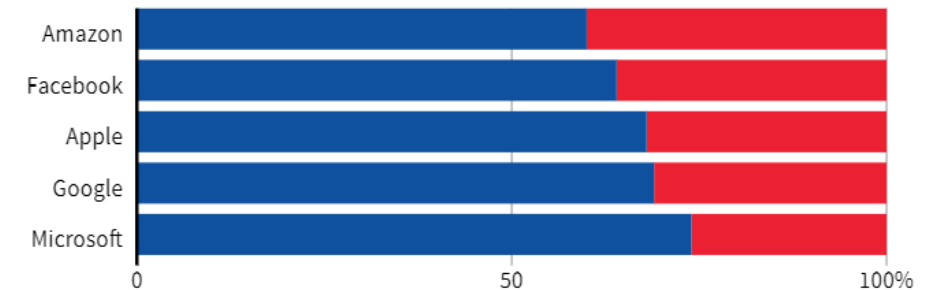
Amazon edited the programs to make them neutral to these particular terms. But that was no guarantee that the machines would not devise other ways of sorting candidates that could prove discriminatory

[After re-training...] the technology favored candidates who described themselves using verbs more commonly found on male engineers' resumes, such as "executed" and "captured," one person said.

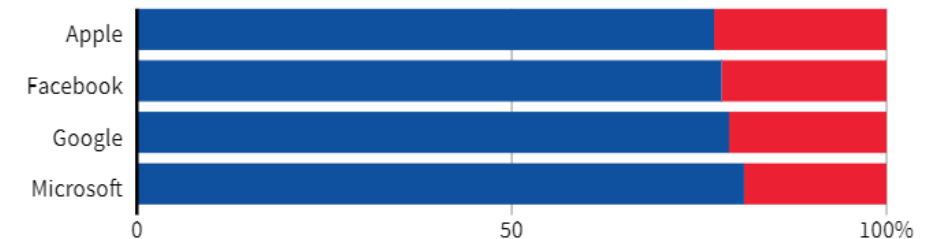
This is an open question in a field called Algorithmic Fairness.

GLOBAL HEADCOUNT

■ Male ■ Female



EMPLOYEES IN TECHNICAL ROLES



<https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scrap-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G>

The Federalist Papers

Probabilistic text analysis

Ignoring the order of words...

What is the probability of any given word that you write in English?

- $P(\text{word} = \text{"the"}) > P(\text{word} = \text{"pokemon"})$
- $P(\text{word} = \text{"Stanford"}) > P(\text{word} = \text{"Cal"})$

Probabilities of *counts* of words = Multinomial distribution



A document is a large multinomial.

(according to the Global Language Monitor, there are 988,968 words in the English language used on the internet.)

Probabilistic text analysis

Probabilities of *counts* of words = Multinomial distribution

Example document:

#words: $n = 48$

“When my late husband was alive he deposited some amount of Money with china Bank in which the amount will be declared to you once you respond to this message indicating your interest in helping to receive the fund and use it for Gods work as my wish.”

$$P \left(\begin{array}{l} \text{bank} = 1 \\ \text{fund} = 1 \\ \text{money} = 1 \\ \text{wish} = 1 \\ \dots \\ \text{to} = 3 \end{array} \middle| \text{spam} \right) = \frac{n!}{1! 1! 1! 1! \dots 3!} p_{\text{bank}}^1 p_{\text{fund}}^1 \dots p_{\text{to}}^3$$

Note: $P(\text{bank} | \text{spam}) \gg P(\text{bank} | \text{writer=you})$

Probabilistic text analysis

Probabilities of *counts* of words = Multinomial distribution

What about probability of those same words in someone else's writing?

- $P(\text{word} = \text{"probability"} \mid \text{writer} = \text{you}) > P(\text{word} = \text{"probability"} \mid \text{non-CS109 student})$

To determine authorship:

1. Estimate $P(\text{word} \mid \text{writer})$ from known writings
2. Use Bayes' Theorem to determine $P(\text{writer} \mid \text{document})$ for a new writing!



Who wrote the Federalist Papers?

See recordings
10e_all...

Up next:
Independent
RVs and Sums!