

Part 1: Theory

L02: Basics of Probability (1/2)

[Random experiment, independence, conditional probability, chain rule, Bayes' theorem, random variables]

Javed Aslam, Wolfgang Gatterbauer

cs7840 Foundations and Applications of Information Theory (fa24)

<https://northeastern-datalab.github.io/cs7840/fa24/>

9/9/2024

Pre-class conversations

- Last class recapitulation
- We strive to keep it interactive, also among us faculty
- Why class slides need up to 2 days after a class
- Office hours: Usually right after class, or via email / Teams
- Organizational matters: Piazza messages? Canvas pictures did not display correctly?
- New class arrivals

- Today:
 - The basics of probability theory

Last time

- Introduction
- Course logistics
- Entropy
 - examples

Today

- Probability primer

Next time

- Basic concepts of Info-Theory

Probability

- Random experiment
- generate outcomes $w \in \Omega$
- sample space: set of all possible outcomes Ω

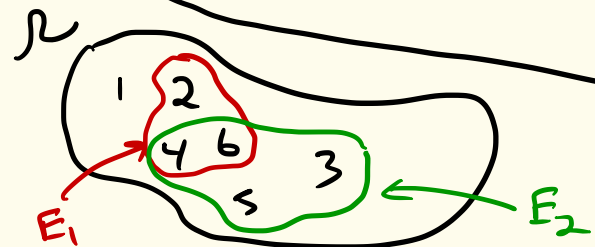
• Event: subset of sample space

• $p: \Omega \rightarrow \mathbb{R}$ probability measure

• $0 \leq p(w) \leq 1 \quad \forall w \in \Omega$

• $\sum_{w \in \Omega} p(w) = 1$

we will assume for now that $p(w) = \frac{1}{|\Omega|} \quad \forall w \in \Omega$



Example

- roll a fair six-sided die
- roll a 5

• $\{1, 2, 3, 4, 5, 6\} = \Omega$

• $E_1 = \text{"even"} = \{2, 4, 6\}$

• $E_2 = \text{"} \geq 3 \text{"} = \{3, 4, 5, 6\}$

• $p(1) = p(2) = p(3) = \dots = p(6) = \frac{1}{6}$

$P(E) = \sum_{w \in E} p(w)$

e.g. $P(E_1) = p(2) + p(4) + p(6) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{1}{2}$

iff $p(w) = \frac{1}{|\Omega|} \quad \forall w \in \Omega$

$P(E) = \frac{|E|}{|\Omega|} \quad P(E_1) = \frac{|E_1|}{|\Omega|} = \frac{3}{6}$

Examples

① Roll one fair die

$$\Omega = \{1, 2, 3, 4, 5, 6\}$$

② Roll two fair die

$$\Omega = \{ (1,1), (1,2), (1,3), \dots, (2,1), (2,2), \dots, (6,6) \}$$

$$= \{1, 2, 3, \dots, 6\} \times \{1, 2, \dots, 6\}$$

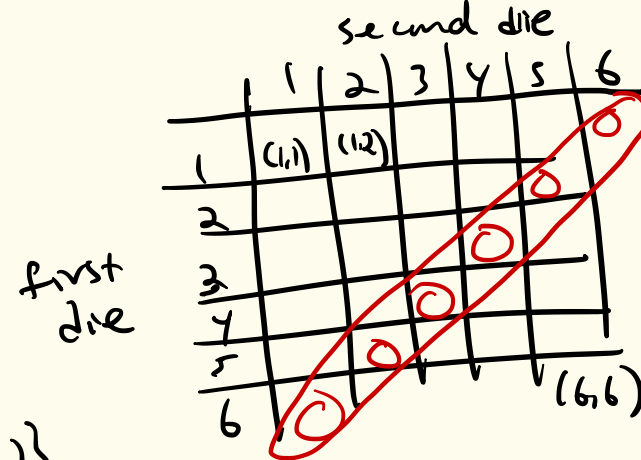
$$E_1 = \text{total is 7} = \{ (1,6), (2,5), (3,4), (4,3), (5,2), (6,1) \}$$

$$P(E_1) = \frac{|E_1|}{|\Omega|} = \frac{6}{36} = \frac{1}{6}$$

$$E_2 = \text{total is greater than 8}$$

$$= 9 \text{ or } 10 \text{ or } 11 \text{ or } 12$$

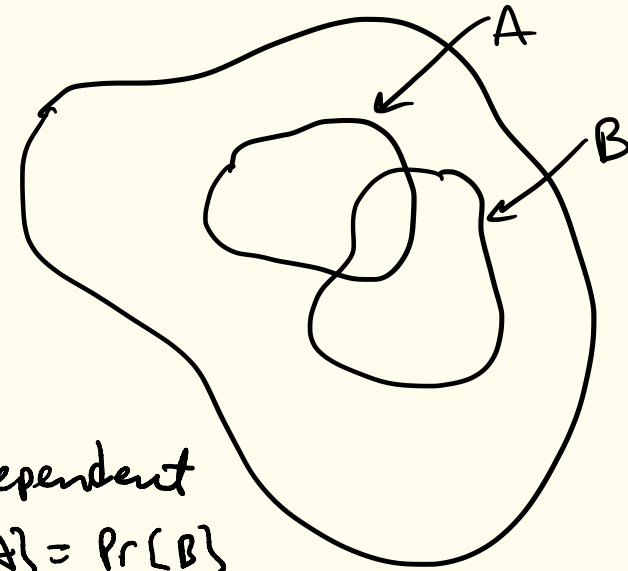
$$P(E_2) = \frac{|E_2|}{|\Omega|} = \frac{10}{36} = \frac{5}{18}$$



$$|E_2| = \begin{array}{cccc} 12 & 11 & 10 & 9 \\ \downarrow & \downarrow & \downarrow & \downarrow \\ 1 & 2 & 3 & 4 \end{array} = 10$$

$$\textcircled{1} \Pr\{A|B\} = \frac{\Pr\{A \cap B\}}{\Pr\{B\}}$$

$$\textcircled{2} \Pr\{B|A\} = \frac{\Pr\{A \cap B\}}{\Pr\{A\}}$$



Independence: Two events A & B are independent if $\Pr\{A|B\} = \Pr\{A\}$, $\Pr\{B|A\} = \Pr\{B\}$

\Rightarrow Knowing B does not change your belief in A .
Knowing A does not change your belief in B .

Claim: A & B are independent if and only if $\Pr\{A \cap B\} = \Pr\{A\} \cdot \Pr\{B\}$

$$\textcircled{1} \Pr\{A|B\} = \frac{\Pr\{A \cap B\}}{\Pr\{B\}}$$

Two events A & B are independent if $\Pr\{A|B\} = \Pr\{A\}$, $\Pr\{B|A\} = \Pr\{B\}$

$$\textcircled{2} \Pr\{B|A\} = \frac{\Pr\{A \cap B\}}{\Pr\{A\}}$$

Claim: A & B are independent if and only if $\Pr\{A \cap B\} = \Pr\{A\} \cdot \Pr\{B\}$

Proof:

(\Rightarrow) If $\Pr\{A|B\} = \Pr\{A\}$, then by $\textcircled{1}$ $\Pr\{A \cap B\} = \Pr\{A|B\} \cdot \Pr\{B\} = \Pr\{A\} \cdot \Pr\{B\}$

(\Leftarrow) If $\Pr\{A \cap B\} = \Pr\{A\} \cdot \Pr\{B\}$, then by $\textcircled{1}$ & $\textcircled{2}$

$$\Pr\{A|B\} = \frac{\Pr\{A \cap B\}}{\Pr\{B\}} = \frac{\Pr\{A\} \cdot \cancel{\Pr\{B\}}}{\cancel{\Pr\{B\}}} = \Pr\{A\}$$

$$\Pr\{B|A\} = \frac{\Pr\{A \cap B\}}{\Pr\{A\}} = \frac{\cancel{\Pr\{A\}} \cdot \Pr\{B\}}{\cancel{\Pr\{A\}}} = \Pr\{B\} \quad \therefore$$

Chain Rules

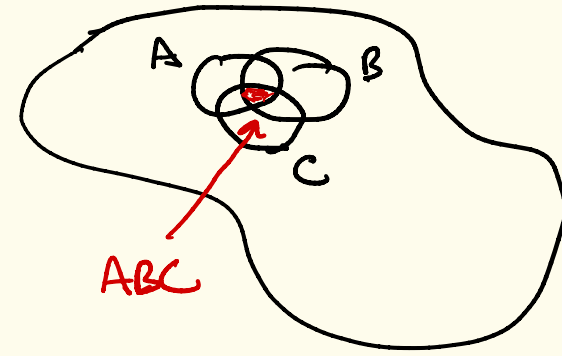
$$\cdot \Pr\{ABC\} = ?$$

- treat AB as an event

$$\Pr\{ABC\} = \Pr\{AB\} \cdot \Pr\{C|AB\}$$

$$= \Pr\{A\} \cdot \Pr\{B|A\} \cdot \Pr\{C|AB\}$$

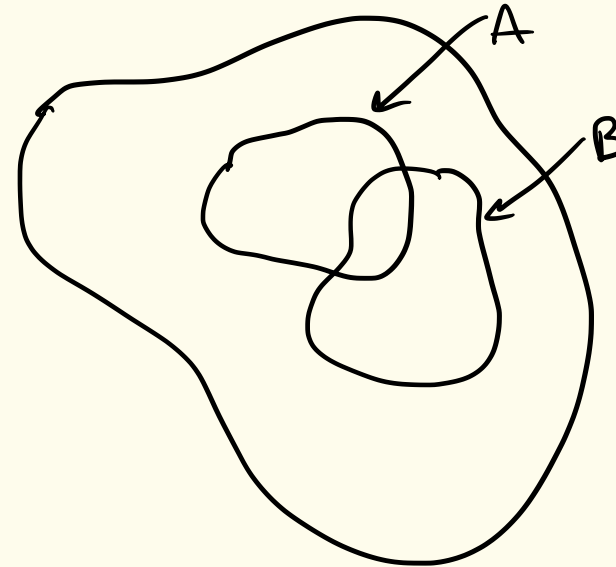
$$\text{Generally: } \Pr\{x_1, x_2, \dots, x_n\} = \Pr\{x_1\} \cdot \Pr\{x_2|x_1\} \cdot \Pr\{x_3|x_1, x_2\} \cdots \Pr\{x_n|x_1, \dots, x_{n-1}\}$$



Bayes Law

$$\textcircled{1} \Pr\{A|B\} = \frac{\Pr\{A \cap B\}}{\Pr\{B\}}$$

$$\textcircled{2} \Pr\{B|A\} = \frac{\Pr\{A \cap B\}}{\Pr\{A\}}$$



$$\Pr\{A|B\} \cdot \Pr\{B\} = \Pr\{A \cap B\} = \Pr\{B|A\} \cdot \Pr\{A\}$$

The equation above shows the relationship between the conditional probabilities and the joint probability. Red curly braces are drawn under the terms $\Pr\{A|B\}$ and $\Pr\{B|A\}$, with circled numbers 1 and 2 respectively above them.

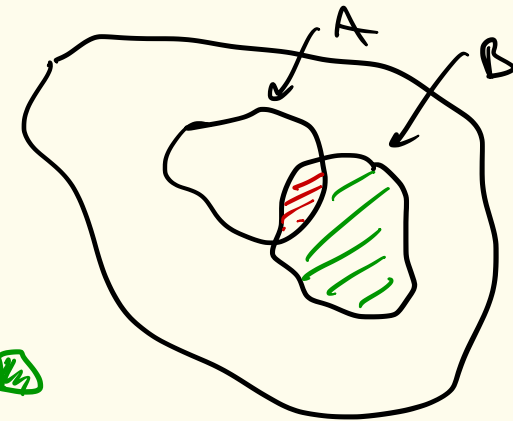
Solving for $\Pr\{A|B\}$...

$$\Pr\{A|B\} = \frac{\Pr\{B|A\} \cdot \Pr\{A\}}{\Pr\{B\}}$$

Bayes Law

$$\Pr\{A|B\} = \frac{\Pr\{B|A\} \cdot \Pr\{A\}}{\Pr\{B\}}$$

what is $\Pr\{B\}$?



$$\Pr\{B\} = \text{red shaded area} + \text{green shaded area}$$

$$= \Pr\{A \cap B\} + \Pr\{\bar{A} \cap B\}$$

$$= \Pr\{B|A\} \cdot \Pr\{A\} + \Pr\{B|\bar{A}\} \cdot \Pr\{\bar{A}\}$$

$$\Pr\{A|B\} = \frac{\Pr\{B|A\} \cdot \Pr\{A\}}{\Pr\{B|A\} \cdot \Pr\{A\} + \Pr\{B|\bar{A}\} \cdot \Pr\{\bar{A}\}}$$

$$\Pr\{H|E\} = \frac{\Pr\{E|H\} \cdot \Pr\{H\}}{\Pr\{E|H\} \cdot \Pr\{H\} + \Pr\{E|\bar{H}\} \cdot \Pr\{\bar{H}\}}$$

H: hypothesis -
patient has Zika

E: evidence -
patient tested
positive on
Zika blood
test

Example: Zika in FL in 2016

- prevalence of Zika in FL $\Pr(\text{Zika}) = 10^{-5}$ (1 in 100,000)
- accuracy of blood test is 99%
 - e.g. $\Pr(\text{pos. test} | \text{Zika}) = 0.99$ or $\Pr(\text{neg. test} | \text{Zika}) = 0.01$
 - $\Pr(\text{pos. test} | \text{no Zika}) = 0.01$ $\Pr(\text{neg. test} | \text{no Zika}) = 0.99$

↙ false negative rate

↘ false positive rate

- patient tests positive: what is chance they have Zika?

$$\Rightarrow \underline{\underline{\text{Not}}} \Pr(\text{pos test} | \text{Zika}) = 0.99$$

$$\Rightarrow \text{You want } \Pr(\text{Zika} | \text{pos. test})$$

$$\Pr\{H|E\} = \frac{\Pr\{E|H\} \cdot \Pr\{H\}}{\Pr\{E|H\} \cdot \Pr\{H\} + \Pr\{E|\bar{H}\} \cdot \Pr\{\bar{H}\}}$$

$$\Pr\{zika | \text{pos test}\} = \frac{\Pr\{\text{pos. test} | zika\} \cdot \Pr\{zika\}}{\Pr\{\text{pos. test} | zika\} \cdot \Pr\{zika\} + \Pr\{\text{pos. test} | \text{not zika}\} \cdot \Pr\{\text{not zika}\}}$$

$$= \frac{0.99 \cdot 10^{-5}}{0.99 \cdot 10^{-5} + 0.01 \cdot (1 - 10^{-5})}$$

$$= \frac{0.0000099}{0.0000099 + 0.0099999}$$

$$\approx 0.00099$$

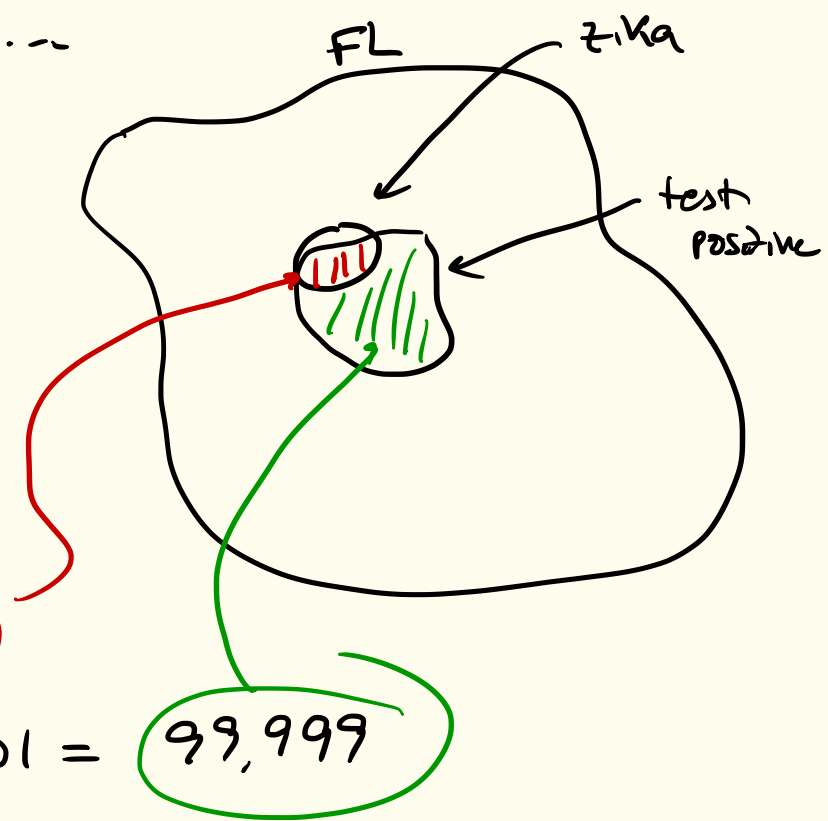
$$\approx 0.1\% \quad \text{i.e. about 1 in 1,000!}$$

Seems wildly counter intuitive, but...

• 10,000,000 people in FL 10^7

• w/ zika? $10^7 \cdot 10^{-5} = 10^2 = 100$

• w/o zika $10^7 - 10^2 = 9,999,900$



test pos w/zika: $100 \cdot 0.99 = 99$

test pos. w/o zika: $9,999,900 \cdot 0.01 = 99,999$

∴ Among those who test pos,
only 99 out of $(99 + 99,999)$
actually have zika - about 1 in 1000.

Random Variable

$$X: \Omega \rightarrow \mathbb{R}$$

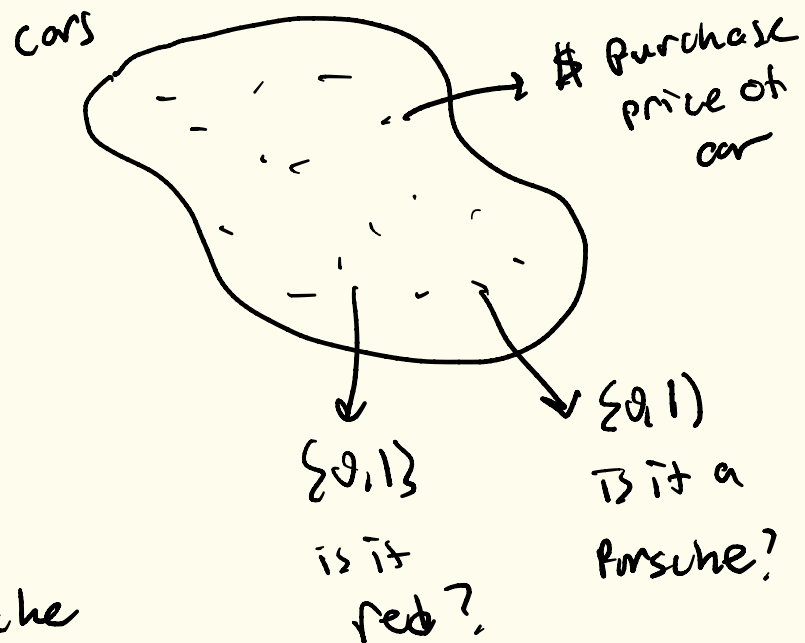
Two primary types:

① "numerical"

e.g. $X(\omega) = \text{price of } \omega \text{ (dollars)}$

② "indicator"

e.g. $X(\omega) = \begin{cases} 1 & \text{if } \omega \text{ is a Porsche} \\ 0 & \text{if not} \end{cases}$



\Rightarrow X is a "random" variable because it depends on the outcome of a random experiment

\Rightarrow Underlying probability measure $p: \Omega \rightarrow \mathbb{R}$ induces a distribution D over the range of the random variable

$$D: \mathbb{R} \rightarrow \mathbb{R}$$

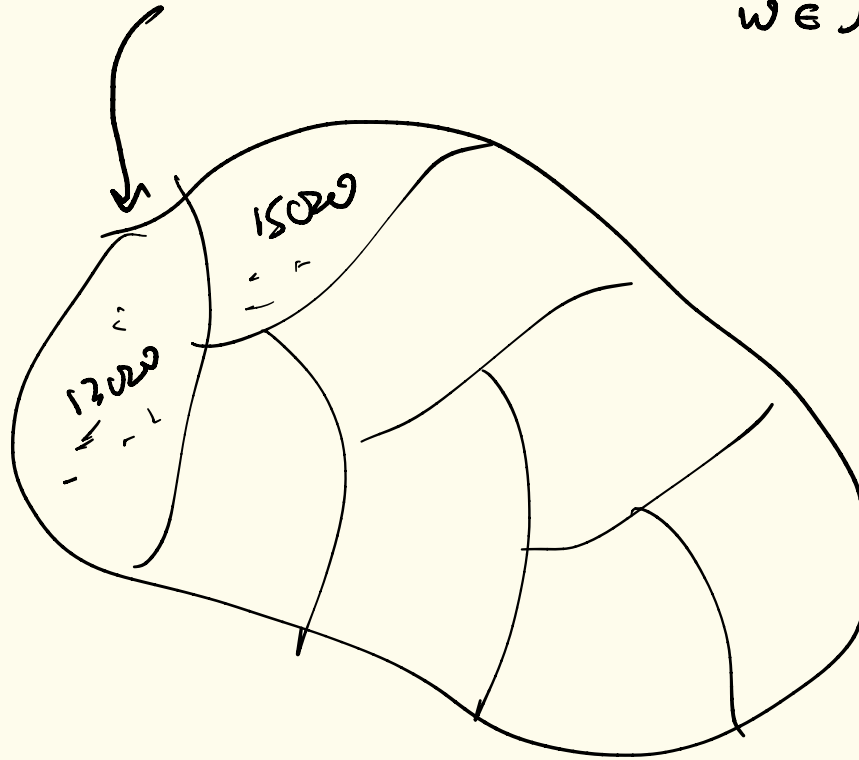
Expectation

- Random variable

$$X: \Omega \rightarrow \mathbb{R}$$

- $E[X] = \sum_x x \cdot \Pr[X=x]$

$$E[X] = \sum_{\omega \in \Omega} X(\omega) \cdot p(\omega)$$



Part 1: Theory

L03: Basics of Probability (2/2)

[Expectation, Variance, Markov chains]

Javed Aslam, Wolfgang Gatterbauer

cs7840 Foundations and Applications of Information Theory (fa24)

<https://northeastern-datalab.github.io/cs7840/fa24/>

9/11/2024

Pre-class conversations

- Last class recapitulation
- PDF Class slides organized by topic (with subsections for classes)
- Any organizational matters: Piazza messages? Organizational matters?
- New class arrivals?

- Today:
 - The basics of probability theory
 - Intuition behind "information" (and "information measures")

Expectation

- Random variable

$$X: \Omega \rightarrow \mathbb{R}$$

$$E\{x\} = \sum_x x \cdot \Pr[X=x]$$

↑

expected or

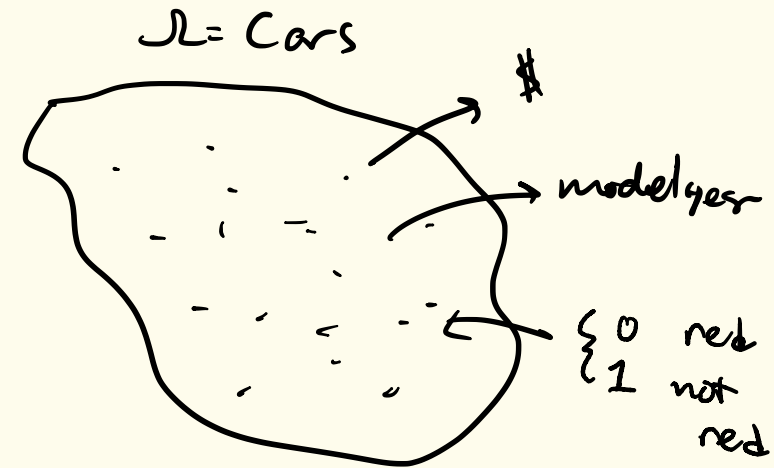
"average" value

$$\frac{1}{2} \text{ cars} \quad \$27,000$$

$$\frac{1}{3} \text{ cars} \quad \$10,000$$

$$\frac{1}{6} \text{ cars} \quad \$15,000$$

$$E\{x\} = \sum_x x \cdot \Pr[X=x] = \frac{1}{2} \cdot 27000 + \frac{1}{3} \cdot 10,000 + \frac{1}{6} \cdot 15000 = \$19,333.33$$



$$E\{x\} = \sum_{\omega \in \Omega} X(\omega) \cdot p(\omega)$$

Example

- Pay \$6 to play game
- Roll two fair 6-sided die
- Pay you sum of die faces except if double, then 0
- $X =$ winnings (profit)

$$E[X] = \sum_x x \cdot \Pr\{X=x\}$$

$$x = -6 \quad \Pr\{X=-6\} = 6/36 = 1/6$$

$$x = -3 \quad \Pr\{X=-3\} = 2/36 = 1/18$$

⋮

$$x = +5 \quad \Pr\{X=5\} = 2/36 = 1/18$$

$$E[X] = \sum_x x \cdot \Pr\{X=x\}$$

$$= (-6) \cdot 1/6 + (-3) \cdot 1/18 + \dots + 5(1/18) = -0.16\bar{6}$$

die 1

	die 2					
	1	2	3	4	5	6
1	-6	-3	-2	-1	0	1
2	-3	-6	-1	0	1	2
3	-2	-1	-6	1	2	3
4	-1	0	1	-6	3	4
5	0	1	2	3	-6	5
6	1	2	3	4	5	-6

$$E[X] = \sum_{\omega \in \Omega} X(\omega) \cdot p(\omega)$$

$$= \frac{1}{36} \sum_{\omega \in \Omega} X(\omega)$$

$$= \frac{1}{36} \left(\text{"sum whole table"} \right)$$

$$= -6/36 = -0.16\bar{6}$$

\Rightarrow lose 16¢
per play, on average

More on Expectation

Example: Roll two fair 6-sided die

Let $X = \text{sum of die faces}$

Q: $E\{X}$?

$$X: \Omega \rightarrow \mathbb{R}$$

e.g. $X(1,2,4) \rightarrow 6$

		die 2					
		1	2	3	4	5	6
die 1	1	(1,1)	(1,2)				
	2	(2,1)			(2,4)		
	3						
	4						
	5						
	6						(6,6)

$$E\{X\} = \sum_x x \cdot \Pr\{X=x\}$$

- $x=2 \quad \Pr\{X=2\} = 1/36$
- $x=3 \quad \Pr\{X=3\} = 2/36$
- $x=4 \quad \Pr\{X=4\} = 3/36$
- \vdots
- $x=12 \quad \Pr\{X=12\} = 1/36$

$$E\{X\} = \sum_x x \cdot \Pr\{X=x\} = 2 \cdot 1/36 + 3 \cdot 2/36 + \dots + 12 \cdot 1/36 = 7$$

$$\begin{aligned} E\{X\} &= \sum_{\omega \in \Omega} X(\omega) \cdot p(\omega) \\ &= \sum_{\omega \in \Omega} X(\omega) \cdot 1/36 \\ &= 1/36 \cdot \sum_{\omega \in \Omega} X(\omega) \\ &= 1/36 \cdot \{\text{sum of table}\} \\ &= 7 \end{aligned}$$

Linearity of Expectation
 Let $X_1 =$ r.v. for first die roll
 Let $X_2 =$ r.v. for second die roll
 Let $X = X_1 + X_2$
 $E\{X\} = E\{X_1 + X_2\} = E\{X_1\} + E\{X_2\} = 3.5 + 3.5 = 7$

Variance & standard deviation

Case 1

$$4'10'' \quad 5' \quad 5'2''$$

$$E\{x_1\} = 5'$$

Case 2

$$4' \quad 5' \quad 6'$$

$$E\{x_2\} = 5'$$

Case 3

$$3' \quad 5' \quad 7'$$

$$E\{x_3\} = 5'$$

Case 2 How to measure "variability"

$$\textcircled{1} E\{Y_1\} = \sum_{\omega \in \Omega} Y(\omega) \cdot p(\omega)$$

$$= -12'' \cdot 1/3 + 0'' \cdot 1/3 + (+12'') \cdot 1/3 \\ = 0''$$

$$\textcircled{2} E\{Y_2\} = |-12''| \cdot 1/3 + |0''| \cdot 1/3 + |12''| \cdot 1/3 \\ = 12 \cdot 1/3 + 0 \cdot 1/3 + 12 \cdot 1/3 = 8''$$

$$\textcircled{3} E\{Y_3\} = (-12'')^2 \cdot 1/3 + (0'')^2 \cdot 1/3 + (12'')^2 \cdot 1/3 \\ = 96 \text{ in}^2$$

3 ways

$$\textcircled{1} Y_1 = X - E\{X\}$$

~~X~~ nes & pos cancel

$$\textcircled{2} Y_2 = |X - E\{X\}|$$

- mean absolute deviation

$$\textcircled{3} Y_3 = (X - E\{X\})^2$$

- variance

\Rightarrow take square root, get standard deviation $\sqrt{96 \text{ in}^2} = 9.8 \text{ in}$

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

$$\Rightarrow \sigma = \sqrt{\text{Var}(X)} = \sqrt{E\{(X - E\{X\})^2\}} \quad - \text{standard deviation}$$

- back in original
units

Claim: $E\{(X - E\{X\})^2\} = E\{X^2\} - (E\{X\})^2$

example: case 2 $E\{X\} = 5' = 60''$

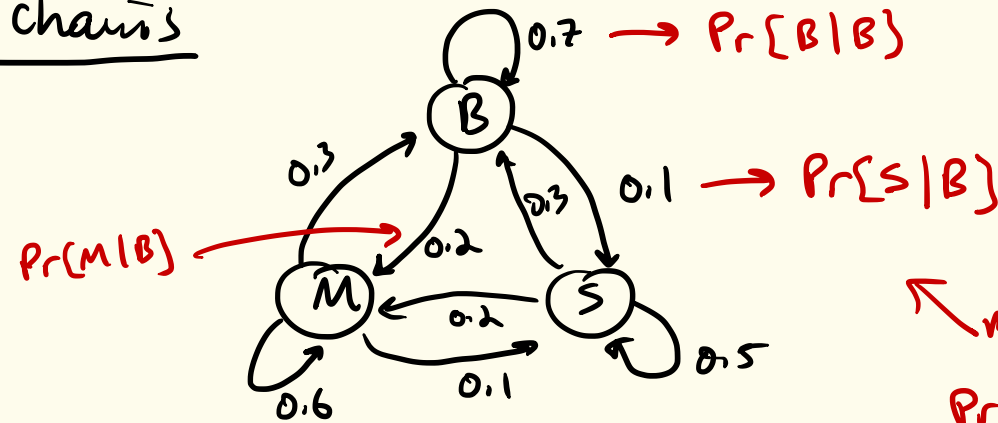
$$E\{X^2\} = \frac{(48'')^2 + (60'')^2 + (72'')^2}{3}$$

$$\text{Claim: } 3696 - 60^2 = 3696 - 3600 = \underline{\underline{96}}$$

$$\begin{aligned} \text{Pf: } E\{(X - E\{X\})^2\} &= E\{X^2 - 2X E\{X\} + (E\{X\})^2\} \\ &= E\{X^2\} - E\{2X E\{X\}\} + E\{(E\{X\})^2\} \\ &= E\{X^2\} - 2E\{X\} \cdot E\{X\} + (E\{X\})^2 \\ &= E\{X^2\} - 2(E\{X\})^2 + (E\{X\})^2 \\ &= E\{X^2\} - (E\{X\})^2 \end{aligned}$$

Markov chains

B: Bertucci's
M: Margaritas
S: Sato



must sum to 1
 $Pr(B|B) + Pr(S|B) + Pr(M|B) = 1$

State transition matrix:

$$P = \begin{matrix} & \begin{matrix} B & M & S \end{matrix} \\ \begin{matrix} B \\ M \\ S \end{matrix} & \begin{pmatrix} .7 & .2 & .1 \\ .3 & .6 & .1 \\ .3 & .2 & .5 \end{pmatrix} \end{matrix} \cdot \begin{matrix} \text{all rows sum to 1} \\ \text{stochastic matrix} \end{matrix}$$

Q: If I were to "run" the Markov Chain, what is the long-term fraction of time I would spend visiting each of the nodes?

A: Stationary distribution $\vec{\pi} = \langle \pi_B, \pi_M, \pi_S \rangle$ $\pi_B + \pi_M + \pi_S = 1$

```
[jay@jay-mac-2021 Downloads % ./restaurantMC.pl 10
```

```
BMBBSSSBMB
```

```
Counts: B=5; M=2; S=3
```

```
Prob: B=0.5; M=0.2; S=0.3
```

```
[jay@jay-mac-2021 Downloads % ./restaurantMC.pl 100
```

```
BSBMBSSSMBBBBSSBSSMMMMSSSMBBBBBMMMMBSSSSMBSSMBBBBBMMMMBBBBMBBMMMMBSMMMMBBBBBBBBB
```

```
Counts: B=37; M=46; S=17
```

```
Prob: B=0.37; M=0.46; S=0.17
```

```
[jay@jay-mac-2021 Downloads % ./restaurantMC.pl 1000
```

```
BBSBMMBSSSSMMBBSBBSBMMBSSSSBBBBSSSMBBBBSSBSSBMSMBSMMMMBSSMMMMBBBBBBBBSSBBMMBMMBMMBMMSSB
BSSSSBBBBBMMBMMMMSSSSSBSSSSBMMMMBMMBBBBBMBSSBSSBMBMMBMMBMMBMMMMBBBBBMBSSBBSSMSBBBBBBMMMMMMBMBB
BBBBMMBMBMMMMMMMMBBBBSSSMMSBBBMMMBMBSMMMMBMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMM
MMMSBMMSSMSBBBMBSSBSBBBBSSMBBMSBMMBMMBMMBMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMM
BBSSMMMMBSSBSSBMMMMMMBMMBBBBSSMMBSSSMBMMBBBBBSSMMBBBBBBBBBBBBBMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMM
SBSSMMMMBSSSMSSSSBBBBMMBSSSSMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMM
BBBBBBBBSSMBBBBBBBBBBSSBSSBSSSSBMMBMMBSSSMBMMBBBBBSSSMBBBBBBBBBBBBBMMMMMMMMMMMMMMMMMMMMMMMM
MMBMBSSBBBBMMBMMBSSBSSBMMBSSSMSSMMMMMMMMSSBBMMBMMBSSMBMMBBBBBSSBBBBMMBMMBMMBMMMMMMMMMMMMMMMM
BBBBBBBBSSSMSSBMMBMMBMMSSSSMMMMSSSSMBMMBBBBSSMMBMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMMM
```

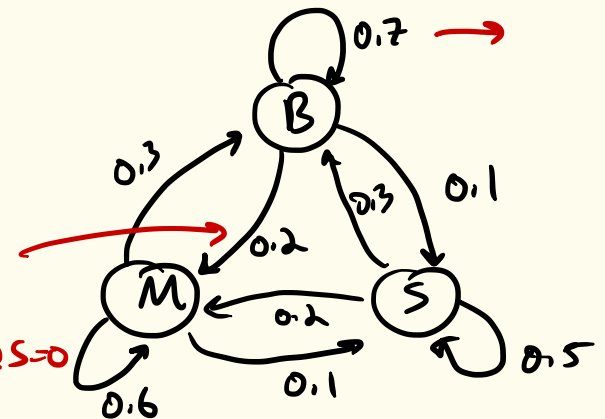
```
Counts: B=488; M=366; S=146
```

```
Prob: B=0.488; M=0.366; S=0.146
```

⇒ way too long to converge, and just for 3 nodes!

① simulate - takes too long to converge

② solve algebraically



① $B = .7 \cdot B + .3 \cdot M + .3 \cdot S$

② $M = 0.2 \cdot B + 0.6 M + 0.2 \cdot S \rightarrow -.2B + .4M - 0.2S = 0$

③ $S = 0.1 \cdot B + 0.1 \cdot M + 0.5 \cdot S \rightarrow -.1B - .1M + .5S = 0$

① $.3B - .3M - .3S = 0$

But also have

② + ③ $.3B - .3M - .3S = 0$

~~④~~ $B + M + S = 1$

$B + M + S = 1$

5x② $-B + 2M - S = 0$

10x③ $-B - M + 5S = 0$

$3M = 1 \Rightarrow$

$M = 1/3$

$6S = 1 \Rightarrow$

$S = 1/6$

$B = 1/2$

Works, but in general, solving n equations in n unknowns takes $O(n^3)$ time, which is infeasible for large n .

Guess an answer for stationary distribution

$$B_0 = \frac{1}{3}$$

$$M_0 = \frac{1}{3}$$

$$S_0 = \frac{1}{3}$$

Try it

$$B_1 = .7 \times \frac{1}{3} + .3 \times \frac{1}{3} + .3 \times \frac{1}{3} = .433$$

$$M_1 = .2 \times \frac{1}{3} + .6 \times \frac{1}{3} + .2 \times \frac{1}{3} = .333$$

$$S_1 = .1 \times \frac{1}{3} + .1 \times \frac{1}{3} + .5 \times \frac{1}{3} = .233$$

new guess

①

$\frac{1}{2}$

②

$\frac{1}{3}$

③

$\frac{1}{6}$

$$B_2 = .7 \times .433 + .3 \times .333 + .3 \times .233 = .473\bar{3}$$

$$M_2 = \dots$$

$$= .33\bar{3}$$

$$S_2 = \dots$$

$$= .193\bar{3}$$

new guess

\Rightarrow iterate

until convergence

```
[jay@jay-mac-2021 Downloads % ./markovIterate.pl 20 transition.txt
```

```
Processing transition matrix...
```

```
0.3333333333333333 0.3333333333333333 0.3333333333333333
0.4333333333333333 0.3333333333333333 0.2333333333333333
0.4733333333333333 0.3333333333333333 0.1933333333333333
0.4893333333333333 0.3333333333333333 0.1773333333333333
0.4957333333333333 0.3333333333333333 0.1709333333333333
0.4982933333333333 0.3333333333333333 0.1683733333333333
0.4993173333333333 0.3333333333333333 0.1673493333333333
0.4997269333333333 0.3333333333333333 0.1669397333333333
0.4998907733333333 0.3333333333333333 0.1667758933333333
0.4999563093333333 0.3333333333333333 0.1667103573333333
0.4999825237333333 0.3333333333333333 0.1666841429333333
0.4999930094933333 0.3333333333333333 0.1666736571733333
0.4999972037973333 0.3333333333333333 0.1666694628693333
0.4999988815189333 0.3333333333333333 0.1666677851477333
0.4999995526075733 0.3333333333333333 0.1666671140590933
0.4999998210430293 0.3333333333333333 0.1666668456236373
0.4999999284172123 0.3333333333333333 0.1666667382494553
0.4999999713668843 0.3333333333333333 0.1666666952997823
0.4999999885467543 0.3333333333333333 0.1666666781199133
0.4999999954187013 0.3333333333333333 0.1666666712479653
0.4999999981674833 0.3333333333333333 0.1666666684991863
```

⇒ rapid convergence;
very efficient