Annealing is needed after milling to soften the metal.

# CS 311: Reasoning with Knowledge and Probability Theory (Review?)

## Admin

- Will have mancala tournament soon (probably Thursday)
  - If you're taking an extension, make sure to get it in by 1:30pm tomorrow
- Schedule
  - Assignment 3 (paper review) out soon, due next Tuesday, solo assignment
  - Look at written problems 2 by Thursday
  - Assignment 4 out next Tuesday
    - Due before spring break
  - First midterm (take-home) week before spring break

## Human agents

How do humans represent knowledge?
  - ontologies
  - scripts

How do humans reason/make decisions?
  - logic
  - probability
  - utility/cost-benefit
  - two decision systems: intuition/reasoning
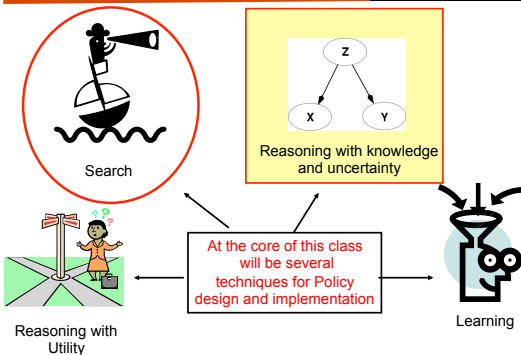    - http://www.princeton.edu/~kahneman/

## An example

Answer the following *as quickly as possible…*

## An example
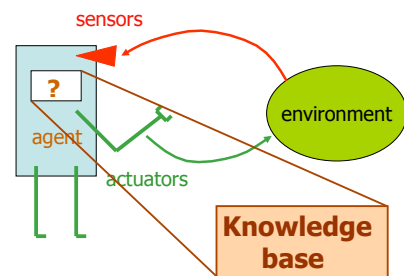
A bat and a ball together cost $1.10. The bat costs a dollar more than the ball. How much does the ball cost?
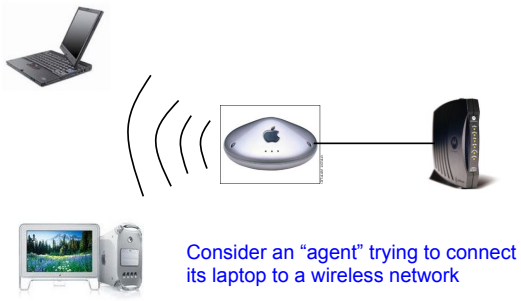
Your first guess is often wrong…

## Policy design: what should an agent do?



Search

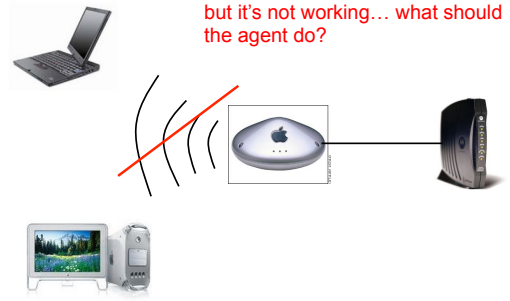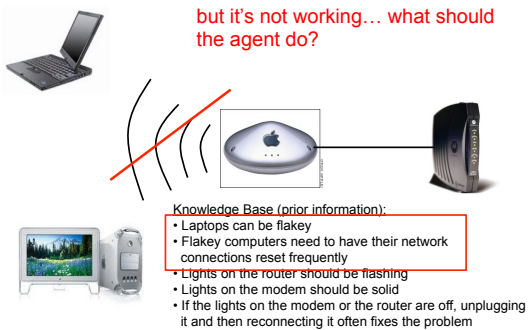Reasoning with knowledge and uncertainty

At the core of this class will be several techniques for Policy design and implementation

Reasoning with Utility

Learning

## Knowledge-Based Agent



sensors

?

agent

actuators

environment

Knowledge base

## Example: Connecting to a home network



Consider an "agent" trying to connect its laptop to a wireless network

## Example: Connecting to a home network



but it's not working… what should the agent do?

## Example: Connecting to a home network



but it's not working… what should the agent do?

Knowledge Base (prior information):
• Laptops can be flakey
• Flakey computers need to have their network connections reset frequently
• Lights on the router should be flashing
• Lights on the modem should be solid
• If the lights on the modem or the router are off, unplugging it and then reconnecting it often fixes the problem

## Example: Connecting to a home network



Knowledge Base (prior information):
• Laptops can be flakey
• Flakey computers need to have their network connections reset frequently
• Lights on the router should be flashing
• Lights on the modem should be solid
• If the lights on the modem or the router are off, unplugging it and then reconnecting it often fixes the problem
• Resetting the computer's network connection did not help
• The lights on the modem are off

## How do we represent knowledge?

Procedurally (HOW):
- Write methods that encode how to handle specific situations in the world
  - `chooseMoveMancala()`
  - `driveOnHighway()`

Declaratively (WHAT):
- Specify facts about the world
  - Two adjacent regions must have different colors
  - If the lights on the modem are off, it is not sending a signal
- Key is then how do we reason about these facts

## Logic for Knowledge Representation

Logic is a declarative language to:

Assert sentences representing facts that hold in a world W (these sentences are given the value true)

Deduce the true/false values to sentences representing other aspects of W

## Propositional logic

Four children have different favorite dinosaurs. Find out who likes which dinosaur.

|  | T. rex | Stegosaurus | Velociraptor | Triceratops |
|---|---|---|---|---|
| Amy |  |  |  |  |
| Bob |  |  |  |  |
| Cal |  |  |  |  |
| Deb |  |  |  |  |

1. Bob's favorite dinosaur does not have an "x" in its name.

2. Amy only likes dinosaurs that walk on four legs.

3. Neither Cal's nor Amy's favorite dinosaur has triangular plates along its back.

4. Bob's favorite dinosaur is a meat-eater.
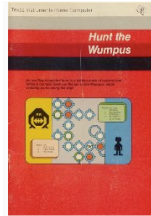
## Propositional logic

Four children have different favorite dinosaurs. Find out who likes which dinosaur.

|  | T. rex | Stegosaurus | Velociraptor | Triceratops |
|---|---|---|---|---|
| Amy |  |  |  |  |
| Bob |  |  |  |  |
| Cal |  |  |  |  |
| Deb |  |  |  |  |

1. Bob's favorite dinosaur does not have an "x" in its name.

2. Amy only likes dinosaurs that walk on four legs.

3. Neither Cal's nor Amy's favorite dinosaur has triangular plates along its back.

4. Bob's favorite dinosaur is a meat-eater.

T.Rex has an x in it
Stegasaurus and Triceritops walk on 4 legs
T.Rex and Velociraptors eat meat
Bob likes Amy
…

## Hunt the Wumpus



Invented in the early 70s (i.e. the "good old days" of computer science)

– originally command-line (think black screen with greenish text)

---

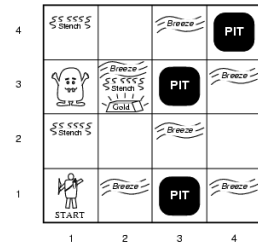## The Wumpus World (as defined by the book)

Performance measure
– gold +1000, death -1000 (falling into pit or eaten by wumpus)
– -1 per step, -10 for using the arrow

Environment
– 4x4 grid of rooms
– Agent starts in [1,1] facing right
– gold/wumpus squares randomly chosen
– Any other room can have a pit (prob = 0.2)

Sensors: Stench, Breeze, Glitter, Bump, Scream

Actuators: Left turn, Right turn, Forward, Grab, Release, Shoot



---

## Wumpus world characterization

### Fully Observable?

– No… until we explore, we don't know things about the world
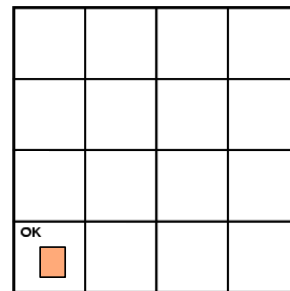
### Deterministic

– Yes

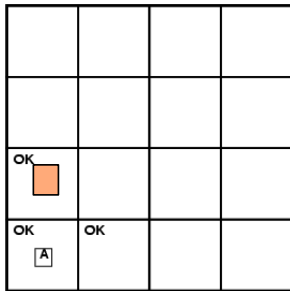### Discrete

– Yes

---

## Exploring a wumpus world



= Agent
B = Breeze
G = Glitter/Gold
OK = Safe Square
P = Pit
S = Stench
W = Wumpus

What do we know?

stench = none, breeze = none, glitter = none, bump = none, scream = none

## Exploring a wumpus world

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| **OK** | | | |
| **OK** A | **OK** | | |

= Agent
B = Breeze
G = Glitter/Gold
OK = Safe Square
P = Pit
S = Stench
W = Wumpus

What do we know?

stench = none, breeze, glitter = none, bump = none, scream = none

## Exploring a wumpus world

| | | | |
|---|---|---|---|
| p? | | | |
| **OK** B | p? | | |
| **OK** A | **OK** | | |

= Agent
B = Breeze
G = Glitter/Gold
OK = Safe Square
P = Pit
S = Stench
W = Wumpus

stench = none, breeze, glitter = none, bump = none, scream = none
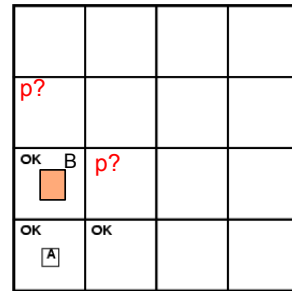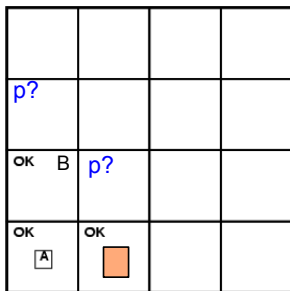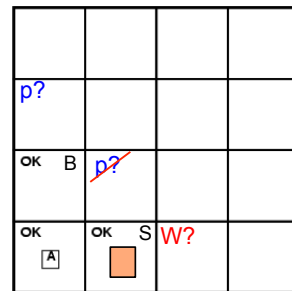
## Exploring a wumpus world

| | | | |
|---|---|---|---|
| p? | | | |
| **OK** B | p? | | |
| **OK** A | **OK** | | |

= Agent
B = Breeze
G = Glitter/Gold
OK = Safe Square
P = Pit
S = Stench
W = Wumpus

What do we know?

stench, breeze = none, glitter = none, bump = none, scream = none

## Exploring a wumpus world

| | | | |
|---|---|---|---|
| p? | | | |
| **OK** B | p? | | |
| **OK** A | **OK** S | W? | |

= Agent
B = Breeze
G = Glitter/Gold
OK = Safe Square
P = Pit
S = Stench
W = Wumpus

stench, breeze = none, glitter = none, bump = none, scream = none

## Exploring a wumpus world

| | | | |
|---|---|---|---|
| | | | |
| p? | | | |
| **OK** B ok | | | |
| | ■ | | |
| **OK** Ⓐ | **OK** S | W? | |

= Agent
B = Breeze
G = Glitter/Gold
OK = Safe Square
P = Pit
S = Stench
W = Wumpus

**What do we know?**

stench = none, breeze = none, glitter = none, bump = none, scream = none

---

## Exploring a wumpus world

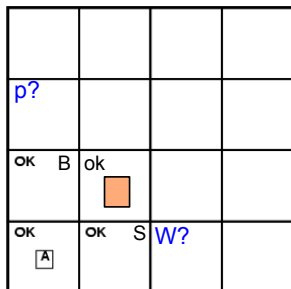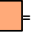| | | | |
|---|---|---|---|
| | | | |
| p? | ok | | |
| **OK** B ok | ok | | |
| | ■ | | |
| **OK** Ⓐ | **OK** S | W? | |

= Agent
B = Breeze
G = Glitter/Gold
OK = Safe Square
P = Pit
S = Stench
W = Wumpus

stench = none, breeze = none, glitter = none, bump = none, scream = none
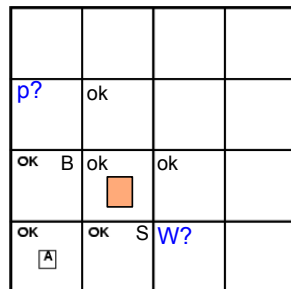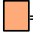
---

## Exploring a wumpus world

| | | | |
|---|---|---|---|
| | | | |
| p? | ok | | |
| **OK** B ok | ok | | |
| | | ■ | |
| **OK** Ⓐ | **OK** S | W? | |

= Agent
B = Breeze
G = Glitter/Gold
OK = Safe Square
P = Pit
S = Stench
W = Wumpus

stench, breeze, glitter, bump = none, scream = none

---

## Wumpus with propositional logic

Using logic statements could determine all of the "safe" squares

### A few problems?

– Sometimes, you have to guess (i.e. no safe squares)
– Sometimes the puzzle isn't solvable (21% of the puzzles are not solvable at all)
– Wumpus may eat you ☺

## Hunt the Wumpus

A modern version…
- http://www.dreamcodex.com/wumpus.php



## Weather rock



## Weather rock

Is the weather always "fine" if the rock is dry?



## Weather rock

## Weather rock

Is the weather always "fine" if the rock is dry
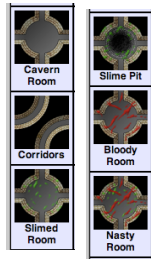AND there isn't an umbrella over it?



## The real world…

Cannot always be explained by rules/facts
– The real world does not conform to logic

Sometimes rocks get wet for other reasons (e.g.
dogs)

Sometimes tomatoes are green, bananas taste like
apples and T.Rex's are vegetarians



## Probability theory

Probability theory enables us to make *rational*
decisions

Allows us to account for uncertainty
– Sometimes rocks get wet for other reasons



## Basic Probability Theory: terminology

An **experiment** has a set of potential outcomes, e.g., throw
a die

The **sample space** of an experiment is the set of all
possible outcomes, e.g., {1, 2, 3, 4, 5, 6}

An **event** is a subset of the sample space.
– {2}
– {3, 6}
– even = {2, 4, 6}
– odd = {1, 3, 5}

We will talk about the probability of events

## Random variables

A random variable is a mapping of all possible outcomes of an experiment to an event

It represents all the possible values of something we want to measure in an experiment

For example, random variable, *X*, could be the number of heads for a coin
- note this is different than the sample space

| space | HHH | HHT | HTH | HTT | THH | THT | TTH | TTT |
|-------|-----|-----|-----|-----|-----|-----|-----|-----|
| *X* | 3 | 2 | 2 | 1 | 2 | 1 | 1 | 0 |

---

## Random variables

We can then talk about the probability of the different values of a random variable

The definition of probabilities over *all* of the possible values of a random variable defines a **probability distribution**

| space | HHH | HHT | HTH | HTT | THH | THT | TTH | TTT |
|-------|-----|-----|-----|-----|-----|-----|-----|-----|
| *X* | 3 | 2 | 2 | 1 | 2 | 1 | 1 | 0 |

| X | P(X) | |
|---|------|---|
| 3 | P(X=3) = | |
| 2 | P(X=2) = | **?** |
| 1 | P(X=1) = | |
| 0 | P(X=0) = | |

---

## Random variables

We can then talk about the probability of the different values of a random variable

The definition of probabilities over *all* of the possible values of a random variable defines a **probability distribution**

| space | HHH | HHT | HTH | HTT | THH | THT | TTH | TTT |
|-------|-----|-----|-----|-----|-----|-----|-----|-----|
| *X* | 3 | 2 | 2 | 1 | 2 | 1 | 1 | 0 |

| X | P(X) |
|---|------|
| 3 | P(X=3) = 1/8 |
| 2 | P(X=2) = 3/8 |
| 1 | P(X=1) = 3/8 |
| 0 | P(X=0) = 1/8 |

---

## Probability distribution

To be explicit:
- A probability distribution assigns probability values to all possible values of a random variable
- These values must be >= 0 and <= 1
- These values must sum to 1 for all possible values of the random variable

Are these probability distributions?

| X | P(X) |
|---|------|
| 3 | P(X=3) = 1/2 |
| 2 | P(X=2) = 1/2 |
| 1 | P(X=1) = 1/2 |
| 0 | P(X=0) = 1/2 |

| X | P(X) |
|---|------|
| 3 | P(X=3) = -1 |
| 2 | P(X=2) = 2 |
| 1 | P(X=1) = 0 |
| 0 | P(X=0) = 0 |

## Unconditional/prior probability

Simplest form of probability is
– P(X)

*Prior probability*: without any additional information, what is the probability
- What is the probability of a heads?
- What is the probability it will rain today?
- What is the probability a student will get an A in AI?
- What is the probability a person is male?
- …

## Joint distributions

We can also talk about probability distributions over multiple variables, called a *joint distribution*

P(X,Y)
- probability of X *and* Y
- a distribution over the cross product of possible values

| AIPass | P(AIPass) |
|--------|-----------|
| true   | 0.89      |
| false  | 0.11      |

| EngPass | P(EngPass) |
|---------|------------|
| true    | 0.92       |
| false   | 0.08       |

| AIPass AND EngPass | P(AIPass, EngPass) |
|--------------------|---------------------|
| true, true         | .88                 |
| true, false        | .01                 |
| false, true        | .04                 |
| false, false       | .07                 |

## Joint distribution

Still a probability distribution
- all values between 0 and 1, inclusive
- all values sum to 1

*All* questions/probabilities of the two variables can be calculate from the joint distribution
- P(X), P(Y), …

| AIPass AND EngPass | P(AIPass, EngPass) |
|--------------------|---------------------|
| true, true         | .88                 |
| true, false        | .01                 |
| false, true        | .04                 |
| false, false       | .07                 |

## Conditional probability

As we learn more information about the world, we can update our probability distribution
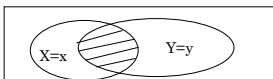- Allows us to incorporate *evidence*

P(X|Y) models this (read "probability of X *given* Y")
- What is the probability of a heads *given* that both sides of the coin are heads?
- What is the probability it will rain today *given* that it is cloudy?
- What is the probability a student will get an A in AI *given* that he/she does all of the written problems?
- What is the probability a person is male *given* that they are over 6 ft. tall?

Notice that the distribution is still over the values of X

## Conditional probability

$$p(X \mid Y) = \frac{P(X,Y)}{P(Y)}$$



Given that Y=y has happened, what proportion of those events does X=x also happen

---

## Conditional probability

$$p(X \mid Y) = \frac{P(X,Y)}{P(Y)}$$



Given that Y=y has happened, what proportion of those events does X=x also happen

| AlPass AND EngPass | P(AlPass, EngPass) |
|---|---|
| true, true | .88 |
| true, false | .01 |
| false, true | .04 |
| false, false | .07 |

What is:
p(AlPass=true | EngPass=false)?

---

## Conditional probability

$$p(X \mid Y) = \frac{P(X,Y)}{P(Y)}$$

| AlPass AND EngPass | P(AlPass, EngPass) |
|---|---|
| true, true | .88 |
| true, false | .01 |
| false, true | .04 |
| false, false | .07 |

What is:
p(AlPass=true | EngPass=false)?

$$\frac{P(true, false) = 0.01}{P(EngPass = false) = 0.01 + 0.07 = 0.08} = 0.125$$

Notice this is different than p(AlPass=true) = 0.89

---

## A note about notation

When talking about a particular assignment, you should technically write p(X=x), etc.

However, when it's clear (like below), we'll often shorten it

Also, we may also say P(X) to generically mean any particular value, i.e. P(X=x)

$$\frac{P(true, false) = 0.01}{P(EngPass = false) = 0.01 + 0.07 = 0.08} = 0.125$$

## Another example

Start with the joint probability distribution:

|  | toothache | | ¬ toothache | |
|---|---|---|---|---|
|  | catch | ¬ catch | catch | ¬ catch |
| cavity | .108 | .012 | .072 | .008 |
| ¬ cavity | .016 | .064 | .144 | .576 |

P(*toothache*) = ?

---

## Another example

Start with the joint probability distribution:

|  | toothache | | ¬ toothache | |
|---|---|---|---|---|
|  | catch | ¬ catch | catch | ¬ catch |
| cavity | .108 | .012 | .072 | .008 |
| ¬ cavity | .016 | .064 | .144 | .576 |

P(*toothache*) = 0.108 + 0.012 + 0.016 + 0.064 = 0.2

---

## Another example

Start with the joint probability distribution:

|  | toothache | | ¬ toothache | |
|---|---|---|---|---|
|  | catch | ¬ catch | catch | ¬ catch |
| cavity | .108 | .012 | .072 | .008 |
| ¬ cavity | .016 | .064 | .144 | .576 |

P(¬*cavity* | *toothache*) = ?

---

## Another example

Start with the joint probability distribution:

|  | toothache | | ¬ toothache | |
|---|---|---|---|---|
|  | catch | ¬ catch | catch | ¬ catch |
| cavity | .108 | .012 | .072 | .008 |
| ¬ cavity | .016 | .064 | .144 | .576 |

$$P(\neg cavity \mid toothache) = \frac{P(\neg cavity, toothache)}{P(toothache)}$$

$$= \frac{0.016 + 0.064}{0.108 + 0.012 + 0.016 + 0.064}$$

$$= 0.4$$

## Normalization

| | toothache | | ¬ toothache | |
|---|---|---|---|---|
| | catch | ¬ catch | catch | ¬ catch |
| cavity | .108 | .012 | .072 | .008 |
| ¬ cavity | .016 | .064 | .144 | .576 |

Denominator can be viewed as a normalization constant α

**P**(*CAVITY* | *toothache*) = α **P**(*CAVITY,toothache*)
= α [**P**(*CAVITY,toothache,catch*) + **P**(*CAVITY,toothache,¬ catch*)]
= α [<0.108,0.016> + <0.012,0.064>]
= α <0.12,0.08> = <0.6,0.4>

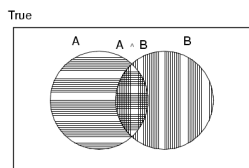unnormalized p(cavity|toothache)      unnormalized p(¬cavity|toothache)

General idea: compute distribution on query variable by fixing evidence variables and summing over hidden/unknown variables

---

## Properties of probabilities

P(*A or B*) = ?

---

## Properties of probabilities

P(*A or B*) = P(*A*) + P(*B*) - P(*A,B*)

True
A    A ∧ B    B

---

## Properties of probabilities

P(¬E) = 1– P(E)

If E1 and E2 are logically equivalent, then:
P(E1)=P(E2).
  – E1: Not all philosophers are more than six feet tall.
  – E2: Some philosopher is not more that six feet tall.
    • Then P(E1)=P(E2).

P(E1, E2) ≤ P(E1).

## The Three-Card Problem

Three cards are in a hat. One is red on both sides (the red-red card). One is white on both sides (the white-white card). One is red on one side and white on the other (the red-white card). A single card is drawn randomly and tossed into the air.

a. What is the probability that the red-red card was drawn?

b. What is the probability that the drawn cards lands with a white side up?

c. What is the probability that the red-red card was not drawn, assuming that the drawn card lands with the a red side up?

## The Three-Card Problem

Three cards are in a hat. One is red on both sides (the red-red card). One is white on both sides (the white-white card). One is red on one side and white on the other (the red-white card). A single card is drawn randomly and tossed into the air.

a. What is the probability that the red-red card was drawn? p(RR) = 1/3

b. What is the probability that the drawn cards lands with a white side up? p(W-up) = 1/2

c. What is the probability that the red-red card was not drawn, assuming that the drawn card lands with the a red side up?

- p(not-RR|R-up)?
- Two approaches:
  - 3 ways that red can be up… of those, only 1 doesn't involve RR = 1/3
  - p(not-RR|R-up) = p(not-RR, R-up) / p(R-up) =
    1/6 / 1/2 = 1/3

## Fair Bets

A bet is fair to an individual I if, according to the individual's probability assessment, the bet will break even in the long run.

Are the following best fair?:

Bet (a):  Win $4.20 if RR;
lose $2.10 otherwise

Bet (b):  Win $2.00 if W-up;
lose $2.00 otherwise

Bet (c):  Win $4.00 if R-up and not-RR;
lose $4.00 if R-up and RR;
neither win nor lose if not-R-up

## Verification

there are six possible outcomes, all equally likely

1. RR drawn, R-up (side 1)
2. RR drawn, R-up (side 2)
3. WR drawn, R-up
4. WR drawn, W-up
5. WW drawn, W-up (side 1)
6. WW drawn, W-up (side 2)

|  | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| a. | $4.20 | $4.20 | -$2.10 | -$2.10 | -$2.10 | -$2.10 |
| b. | -$2.00 | -$2.00 | -$2.00 | $2.00 | $2.00 | $2.00 |
| c. | -$4.00 | -$4.00 | $4.00 | $0.00 | $0.00 | $0.00 |

## Verification

|    | 1        | 2        | 3        | 4        | 5        | 6        |
|----|----------|----------|----------|----------|----------|----------|
| a. | $4.20    | $4.20    | -$2.10   | -$2.10   | -$2.10   | -$2.10   |
| b. | -$2.00   | -$2.00   | -$2.00   | $2.00    | $2.00    | $2.00    |
| c. | -$4.00   | -$4.00   | $4.00    | $0.00    | $0.00    | $0.00    |

expected values:

$$E(a) = \frac{1}{6}4.2 + \frac{1}{6}4.2 + \frac{1}{6}(-2.1) + \frac{1}{6}(-2.1) + \frac{1}{6}(-2.1) + \frac{1}{6}(-2.1) = 0$$

$$E(b) = \frac{1}{6}(-2) + \frac{1}{6}(-2) + \frac{1}{6}(-2) + \frac{1}{6}2 + \frac{1}{6}2 + \frac{1}{6}2 = 0$$

$$E(c) = \frac{1}{6}(-4) + \frac{1}{6}(-4) + \frac{1}{6}4 + \frac{1}{6}0 + \frac{1}{6}0 + \frac{1}{6}0 = -2/3$$

## Why take a bad bet?

| Statistical Edges Against the Player for Casino Games | % |
|---|---|
| Baccarat | 1.17-14.1 |
| Blackjack | |
| Normal | 10.0-20.0 |
| Perfect strategy | 1.2-2.0 |
| Strict card counting | 0.0-2.0 |
| Craps | |
| Normal bets | 1.4-16.7 |
| Single odds | 0.8 |
| Double odds | 0.6 |
| Ten times odds | 0.0 |
| Keno | 29.5 |
| Roulette | 5.26 |
| Slot Machines | 2.0-35.0 |
| Sports Betting | |
| Football and Basketball | 4.54 |
| Single bets | 10.0 |
| Two-bet parlays | 12.5 |
| Three-bet parlays | 31.3 |
| Horse Racing | 19.0 |

http://www.pbs.org/wgbh/pages/frontline/shows/gamble/odds/odds.html

## Monty Hall

- 3 doors
  - behind two, something bad
  - behind one, something good
- You pick one door, but are not shown the contents
- Host opens one of the other two doors that has the bad thing behind it (he always opens one with the bad thing)
- You can now switch your door to the other unopened. Should you?

## Monty Hall

p(win) initially?
  – 3 doors, 1 with a winner, p(win) = 1/3

p(win | shown_other_door)?
  – One reasoning:
    • once you're shown one door, there are just two remaining doors
    • one of which has the winning prize
    • 1/2

  This is not correct!

## Be careful! – Player picks door 1

| | winning location | | host opens | |
|---|---|---|---|---|
| | | 1/2 | Door 2 | |
| 1/3 | Door 1 | 1/2 | Door 3 | |
| | | | | |
| 1/3 | Door 2 | 1 | Door 3 | In these two cases, switching will give you the correct answer. Key: host knows where it is. |
| 1/3 | Door 3 | 1 | Door 2 | |

## Another view

1000 doors
– behind 999, something bad
– behind one, something good

You pick one door, but are not shown the contents

Host opens 998 of the other 999 doors that have the bad thing behind it (he always opens ones with the bad thing)

In essence, you're picking between it being behind your one door or behind any one of the other doors (whether that be 2 or 999)