

## TRICK OR TREAT PROBLEM

AKA "LILY MEETS THE CAT"

140 CONDOS, EACH WITH PORCH LIGHT

6pm - 7pm

UNCERTAINTY:

- SOMEONE AT HOME?
- IF HOME, TIME TO ANSWER?

Noisy indication

RING OUT

ANSWER  $\equiv$  PROVIDE CANDY  $\Rightarrow$  REWARD  $+1$

ACTIONS:

- RING: ONLY ONCE PER HOME [RB]
- WAIT [WH]
- NEXT [NH]

TIME: DISCRETIZED (10sec)

**From:** Pablo Rodriguez Bertorello <pablo@1now.ai>  
**Sent:** Monday, December 4, 2017 8:56 PM  
**To:** usa5105@fedex.com  
**Subject:** ANOTHER

```
[#####
##### TRICK OR TREAT? LILY MEETS CAT
#####
# QUESTION 4
#####

### Given facts
reward_if_answer = 1
probability_being_home = 0.85 # light is on

### Theta
# defined in Q1, maximum likelihood calculation Q2
theta = 0.29

### Probability of answer
# Function of Theta, and iteration number
function probability_of_answer(t)
    if t==0
        return theta
    end
    return (1-theta)^t * theta
end

### Discounted utility
function discounted_utility(utility_table, t, discount_factor)
    if t==1
        return 0 # before first time increment, utility is 0
    end
    return discount_factor * utility_table[t-1]
end

### POLICY ITERATION
function policy_utility(discount_factor, number_of_increments)
    utility_table = fill(0., number_of_increments) # dynamic programming; and track it over time
    for t = 1:number_of_increments
        instant_reward = probability_being_home * probability_of_answer(t) * reward_if_answer
        println("t=", t, " instant_reward: ", instant_reward)
        utility_table[t] = instant_reward + discounted_utility(utility_table, t, discount_factor)
    end
    return utility_table
end

#####
# PART 4a: policy to wait for eternity at house with light (discounted utility)
## Facts
discount_factor = 0.9      # given
number_of_increments = 100  # Run for eternity

## utility over infinite time converges to zero, due to discounting
utility = policy_utility(discount_factor, number_of_increments)

println("UTILITY OVER TIME: ", utility)

## instant reward over time, with decreasing probability of answer
#t=1 instant_reward: 0.2058999999999997
#t=2 instant_reward: 0.1461889999999999
#t=3 instant_reward: 0.1037941899999997
```

PABLO RODRIGUEZ BERTORELLO

2/19

```

#t=4 instant_reward: 0.07369387489999998
#t=5 instant_reward: 0.052322651178999986
#t=6 instant_reward: 0.03714908233708998
#t=7 instant_reward: 0.02637584845933389
#t=8 instant_reward: 0.01872685240612706
#t=9 instant_reward: 0.013296065208350211
#t=10 instant_reward: 0.00944020629792865

# UTILITY OVER TIME (converges to zero): [0.2059, 0.331499, 0.402143, 0.435623, 0.444383, 0.437094, 0.41976, 0.396511, 0.370156, 0.342581, 0.315025,
#0.288282, 0.262832, 0.238948, 0.216756, 0.19629, 0.17752, 0.160377, 0.144772, 0.130602, 0.11776, 0.106139, 0.0956353, 0.0861498,
#0.0775903, 0.0698706, 0.0629115, 0.0566402, 0.0509903, 0.0459012, 0.0413182, 0.0371914, 0.0334759, 0.0301308, 0.0271196, 0.0244089,
#0.0219689, 0.0197727, 0.0177959, 0.0160166, 0.0144152, 0.0129738, 0.0116765, 0.010509, 0.00945814, 0.00851236, 0.00766116, 0.00689506,
#0.00620557, 0.00558503, 0.00502653, 0.00452388, 0.0040715, 0.00366435, 0.00329792, 0.00296813, 0.00267132, 0.00240418, 0.00216377,
#0.00194739, 0.00175265, 0.00157739, 0.00141965, 0.00127768, 0.00114992, 0.00103492, 0.000931431, 0.000838288, 0.000754459, 0.000679013,
#0.000611112, 0.000550001, 0.000495001, 0.000445501, 0.000400951, 0.000360856, 0.00032477, 0.000292293, 0.000263064, 0.000236757,
#0.000213082, 0.000191773, 0.000172596, 0.000155337, 0.000139803, 0.000125823, 0.00011324, 0.000101916, 9.17247e-5, 8.25522e-5,
#7.4297e-5, 6.68673e-5, 6.01805e-5, 5.41625e-5, 4.87462e-5, 4.38716e-5, 3.94845e-5, 3.5536e-5, 3.19824e-5, 2.87842e-5]

##### PART 4b: policy to wait for eternity at house with light (no discounting)
### Facts
discount_factor = 1.0      # given
number_of_increments = 100    # Run for eternity

## utility over infinite time converges to zero, due to discounting
utility = policy_utility(discount_factor, number_of_increments)

println("UTILITY OVER TIME: ", utility)
# UTILITY OVER TIME (converges to >0): [0.2059, 0.352089, 0.455883, 0.529577, 0.5819, 0.619049, 0.645425, 0.664151, 0.677448,
#0.686888, 0.69359, 0.698349, 0.701728, 0.704127, 0.70583, 0.707039, 0.707898, 0.708508, 0.70894, 0.709248, 0.709466, 0.709621,
#0.709731, 0.709809, 0.709864, 0.709904, 0.709932, 0.709951, 0.709966, 0.709976, 0.709983, 0.709988, 0.709991, 0.709994, 0.709996,
#0.709997, 0.709998, 0.709998, 0.709999, 0.709999, 0.709999, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71,
#0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71,
#0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71, 0.71]

##### PART 4c: policy to wait for 1 minute at house with light (no discounting)
### Facts
discount_factor = 1.0      # given
number_of_increments = 6      # Run for eternity

## utility over infinite time converges to zero, due to discounting
utility = policy_utility(discount_factor, number_of_increments)

println("UTILITY OVER TIME: ", utility)
# UTILITY OVER TIME (converges to >0 about same as infinite wait): [0.2059, 0.352089, 0.455883, 0.529577, 0.5819, 0.619049]

#####
# QUESTION 5

##### Part 5a: policy to wait up to eternity at set of houses (no discounting)
### Facts
discount_factor = 1.0      # given
theta = 0.5      # given

number_of_increments = 18    # Run for eternity; if get candy go to next house

## utility over infinite time converges to zero, due to discounting
utility = policy_utility(discount_factor, number_of_increments)

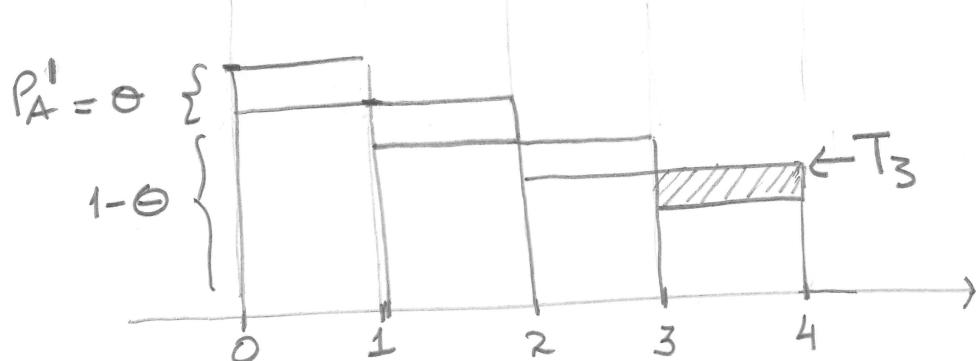
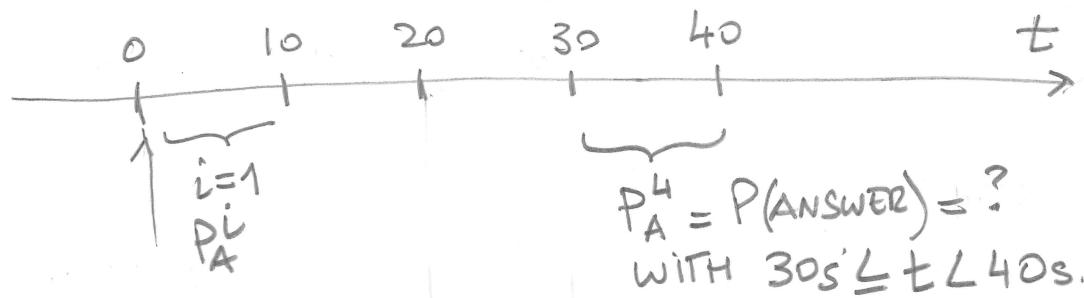
println("UTILITY OVER TIME: ", utility)
# UTILITY OVER TIME (converges to >0 about same as infinite wait): [0.2059, 0.352089, 0.455883, 0.529577, 0.5819, 0.619049]

```

①

$\Theta = P(\text{OCCUPANT ANSWER})$  WITHIN TIME INTERVAL (10s)

✓ INTERVAL (IF NOT OPENED ALREADY)



$T_0 = 1$     ASSUME OCCUPANT KNOWN TO BE HOME

$$T_1 = T_0 \cdot (1-\Theta)$$

$$T_n = T_{n-1} (1-\Theta) = (1-\Theta)^n$$

$$P_A^n = \frac{T_{n-1}}{T_0}, \Theta = \boxed{(1-\Theta)^{n-1}}, \Theta = P_A^n \quad \text{for } n > 1$$

$$\boxed{P_A^4 = (1-\Theta)^3 \cdot \Theta}$$

② TIME TO ANSWER (SECONDS) AT 8 HOMES:

13, 17, 5, 41, 11, 89, 22, 34

INTERVAL	1	2	3	4	5	6	7	8	9
COUNT <sub>i</sub>	1	3	1	1	1	0	0	0	1
OUT OF	8	7	4	3	2	1	1	1	1
$\theta_L^i$	$\frac{1}{8}$	$\frac{3}{7}$	$\frac{1}{4}$	$\frac{1}{3}$	$\frac{1}{2}$	0	0	0	1

$\theta_L^i \equiv P(\text{ANSWER}) \text{ within INTERVAL}$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} P(D|\theta) = \frac{1}{n} \sum_{i=1}^n \theta_L^i$$

$$\hat{\theta} = \frac{1}{9} \cdot \left( \frac{1}{8} + \frac{3}{7} + \frac{1}{4} + \frac{1}{3} + \frac{1}{2} + 1 \right)$$

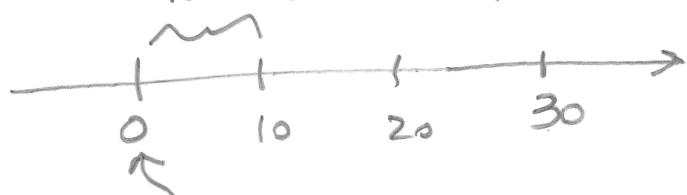
$$\hat{\theta} = \frac{443}{1512} = 0.29$$

## PRIOR PROBABILITY

31

LIGHT	$P(H L)$
ON	0.85
OFF	0.05

NO ANSWER  $\equiv \alpha_1^o$



LIGHT[ON]  $\Rightarrow P(+)=0.85$

$$P(H | a_1^\circ a_2^\circ a_3^\circ) = P(H) = 0.85$$

↑  
CONDITIONAL INDEPENDENCE

3<sup>11</sup>

$$P(\text{ANSWER} | a_4^\circ, a_5^\circ, a_6^\circ) = P_A^7 \cdot \underbrace{P(H|L)}_{0.85 \text{ LIGHT}} \quad \boxed{\text{LIGHTEON}}$$
$$\boxed{P = 0.85 \cdot (1-\theta)^6 \theta}$$

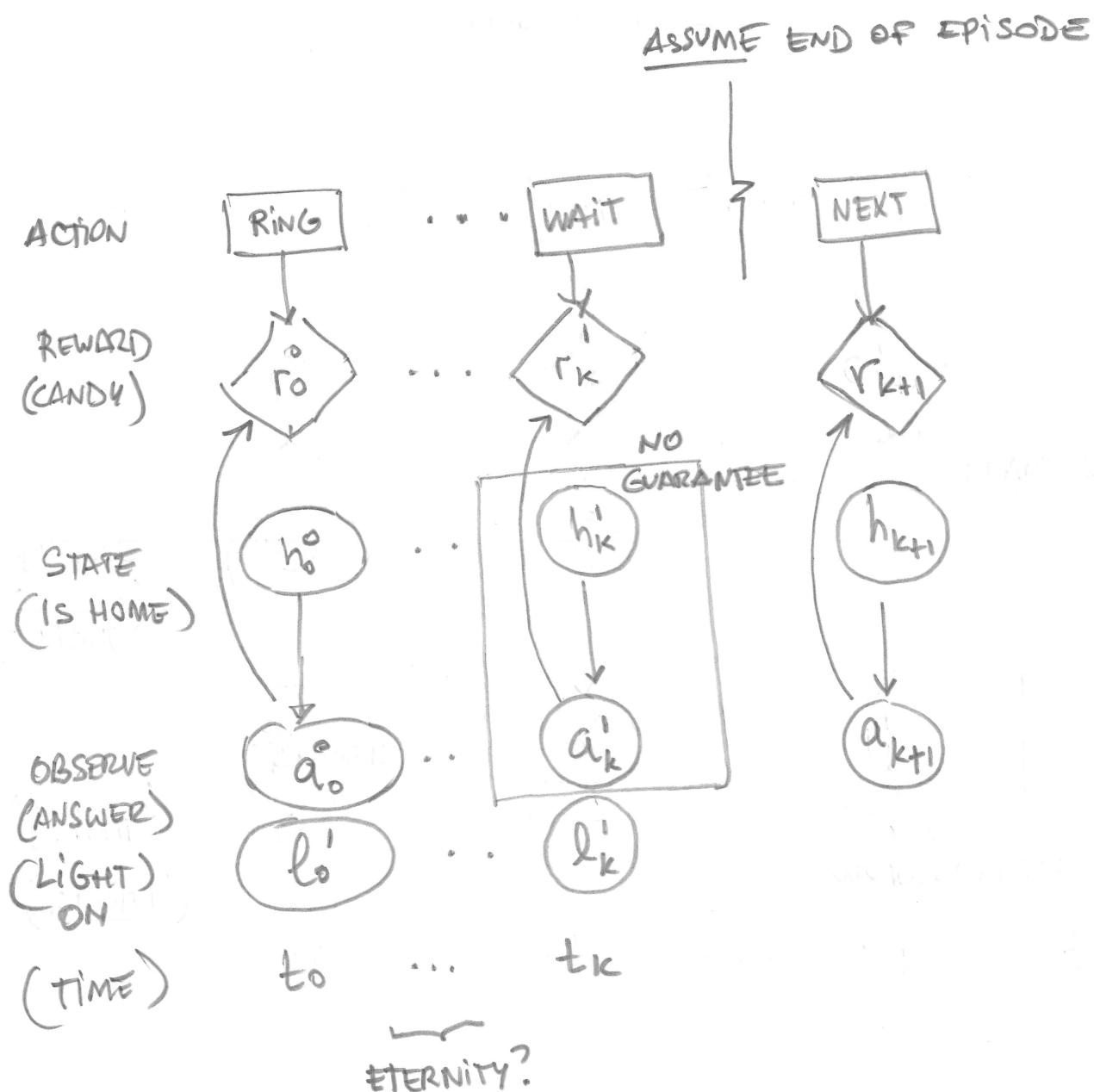
$P_A^7$  PER Q1 ANSWER

④

DOOR W/LIGHT ON

WILLING TO WAIT ETERNITY (NO GUARANTEE)

— EXACT ANSWER IN  $\Theta$  —



$$(4) \quad \text{EU} = ? \quad \gamma = 0.9 \quad \text{LIGHT} = \boxed{\text{ON}} \quad \boxed{1 \equiv (s=s')}$$

$$U_t^{\pi}(s) = R[s, \pi(s)] + \gamma \sum_{s'} T[s'|s, \pi(s)] U_{t+1}^{\pi}(s'),$$

$\underbrace{A = \text{RING/WAIT}}$        $\underbrace{\text{WAIT FOR ETERNITY}}$   
 $\underbrace{0.85(\text{LIGHT})}$

$$R_t[s, \pi(s)] = r(t) \cdot P_A^t \cdot P(H|L)$$

$r=1$        $\underbrace{\text{PROBABILITY OF OCCUPANT}}$   
 $\text{GIVEN}$        $\underbrace{\text{ANSWER IN INTERVAL } t}$   
 $(Q1)$

$$(4a) \quad \boxed{U_t^{\pi}(s) = 0.85 P_A^t + \gamma U_{t-1}^{\pi}(s)} \quad \begin{matrix} \downarrow \\ \text{GIVEN} \end{matrix} \quad U_0^{\pi}(s) = 0$$

WAIT ETERNITY OR

$$\boxed{U_1^{\pi}(s) = 0.85 \cdot \theta}$$

$\underbrace{P_A^1}$

$$U_2^{\pi}(s) = 0.85 \cdot \underbrace{(1-\theta)}_{P_A^{2:1}} \cdot \theta + 0.9 \theta$$

$\underbrace{U_1^{\pi}(s)}$

$$\boxed{U_{\infty}^{\pi}(s) = 0}$$

CODE ATTACHED  
WITH VALUES

$\left\{ \begin{array}{l} \lim_{t \rightarrow \infty} P_A^t = 0 \\ \lim_{t \rightarrow \infty} \gamma^t = 0 \end{array} \right.$

4"

## NO DISCOUNTING

ETERNITY OK  
NO DISCOUNT

(b)

$$U_t^{\pi}(s) = 0.85 P_A^t + U_{t-1}^{\pi}(s) \quad \text{FROM PART(a)}$$

$\uparrow$   
 $\gamma = 1$

$$U_t^{\pi}(s) = 0.85 (1-\theta)^{t-1} \theta + U_{t-1}^{\pi}(s) \quad t > 1$$

FROM Q1

$$U_{\infty}^{\pi}(s) = 0.60$$

> 0 CONVERGENCE  
w/o DISCOUNTING

CODE / VALUES ATTACHED

(c) NO DISCOUNTING

MAX WAIT ONE MINUTE

$$U_6^{\pi}(s) = 0.85 \cdot (1-\theta)^5 \theta + U_5^{\pi}(s)$$

Q4 b

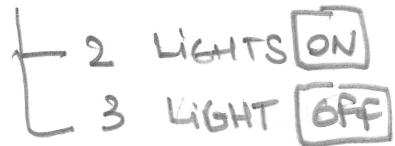
$$U_6^{\pi}(s) = 0.52$$

NEARLY AS MUCH AS 4b  
WITHOUT WAITING ETERNITY

CODE / VALUES ATTACHED.

⑤ IT's 6:57 pm (3 MINUTES REMAIN)

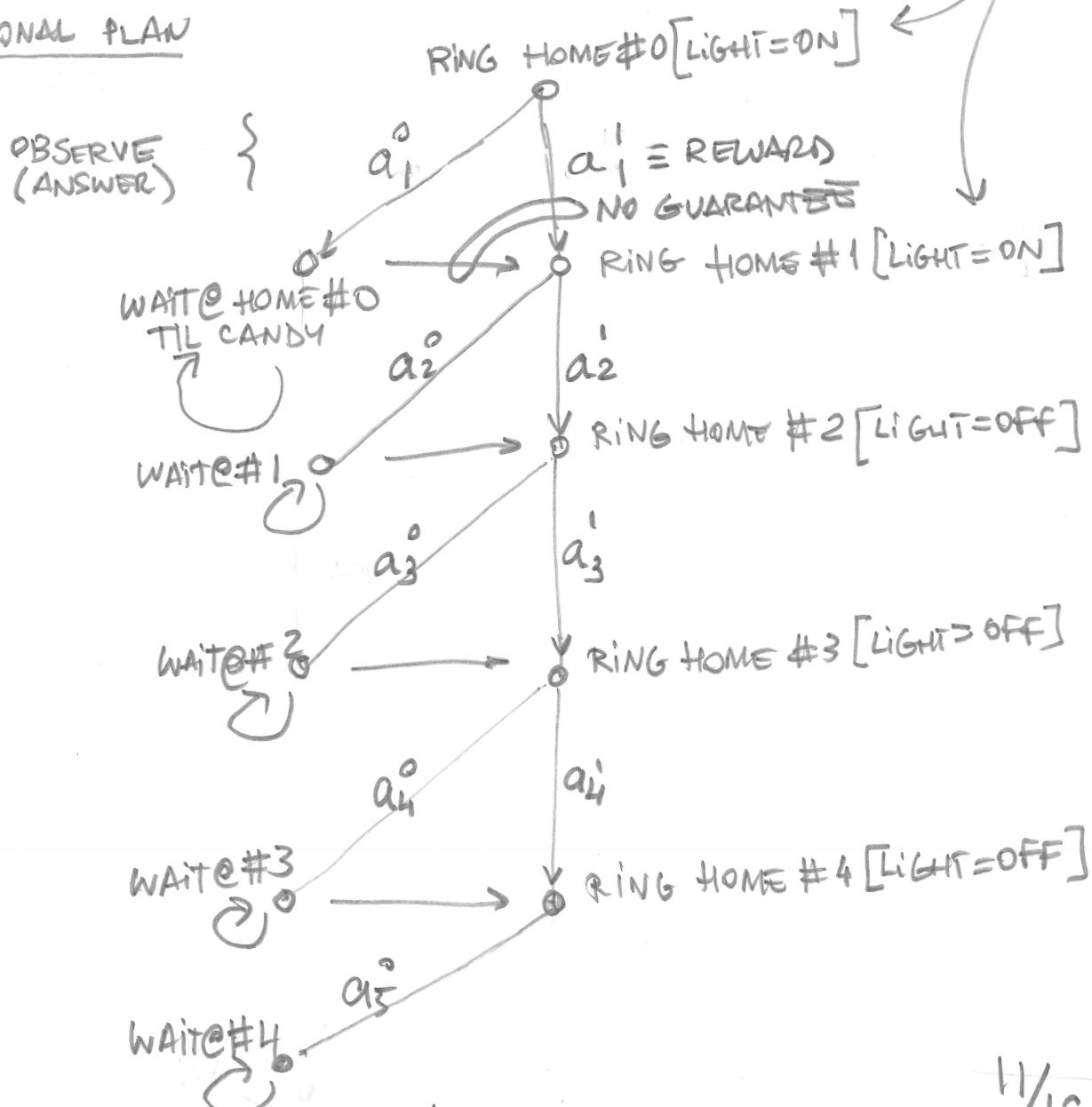
5 HOMES REMAIN



$\Rightarrow \pi \in \text{WAIT TIL GET CANDY. } \mathbb{E} U?$

$\theta = 0.5$   $\gamma = 1$  SMART ORDER  $\equiv$  LIGHT ON FIRST  
— EXACT #, NOT APPROXIMATION —

### CONDITIONAL PLAN



⑤'  $\gamma = 1$   $\theta = 0.5$  WAIT TIL GET CANDY

$$\vec{b}(s) = [P(h^0), P(h^1)]$$

$$\left. \begin{array}{l} b(H | \text{LIGHT=ON}) = [0.15, 0.85] \\ b(H | \text{LIGHT=OFF}) = [0.95, 0.05] \end{array} \right\} \begin{array}{l} P(h^0) \quad \text{GIVEN } P(h^1) \end{array}$$

$$R_t^\pi[s, \text{WAIT}] = 1 \cdot P_A^t$$

$\underbrace{\phantom{0}}_{\text{CANDY}} \quad \underbrace{\phantom{0}}_{\text{ANSWER IN INTERVAL } t}$

$$\left. \begin{array}{l} \alpha_{\text{WAIT}}^{(t)} = [0, (1-\theta)^t \cdot \theta] \\ \alpha_{\text{RING NEXT}}^{(t=1)} = [0, \theta] \end{array} \right\} \begin{array}{l} \text{REWARDS FOR} \\ \text{ACTIONS} \end{array}$$

$\underbrace{\phantom{0}}_{h^0} \quad \underbrace{\phantom{0}}_{h^1}$

5<sup>n</sup>

$U(b)$

A graph showing the relationship between  $U(b)$  (Y-axis) and  $P(\text{HOME})$  (X-axis). The X-axis ranges from  $h^0$  to  $h^1$ , with points  $0.15$ ,  $0.85$  marked. The Y-axis ranges from  $0$  to  $\theta$ , with points  $0.12\theta$ ,  $0.85\theta$  marked. Two curves are shown:  $\alpha(t=1)$  (upper) and  $\alpha(t=2)$  (lower). A vertical double-headed arrow labeled "OVER TIME WAITING" connects the two curves. A box labeled "P(HOME)" is at  $P(\text{HOME}) = 0.15$ . A box labeled "LIGHT = ON" is at  $P(\text{HOME}) = 0.85$ .

$U(t=1)$  OVERWHELMs  $U(t>1)$

$U(b)$

A graph showing the relationship between  $U(b)$  (Y-axis) and  $P(\text{HOME})$  (X-axis). The X-axis ranges from  $h^0$  to  $h^1$ , with points  $0.05$ ,  $0.95$  marked. The Y-axis ranges from  $0$  to  $\theta$ , with point  $0.05\theta$  marked. Two curves are shown:  $\alpha(t=1)$  (upper) and  $\sim \alpha(t=2)$  (lower). A vertical double-headed arrow labeled "OVER TIME WAITING" connects the two curves. A box labeled "P(HOME)" is at  $P(\text{HOME}) = 0.05$ . A box labeled "LIGHT = OFF" is at  $P(\text{HOME}) = 0.95$ .

$U(\text{LIGHT } \boxed{\text{ON}})$  OVERWHELMs  $U(\text{LIGHT } \boxed{\text{OFF}})$

PABLO RODRIGUEZ BERTORELLO

13/19

(5'''')

POLICY: WAIT TIL GET CANDY [MAX 3 MINUTES]

SMART: Go to houses with light on first

$$\theta = 0.5 \quad \gamma = 1$$

$$P(H | \text{LIGHT ON})$$

~ ~

$$U_t^H(s) < 0.85 (1-\theta)^{t-1} \theta + U_{t-1}(s) \quad \text{for } t > 1$$

$$U_0 = 0$$

$$U_1 = 0.85\theta$$

SEE CODE AND VALUES

$$U_{18}(s) \geq 0.42$$

UTILITY MAY BE GREATER  
IF CAN GO TO OTHER HOMES  
(NO GUARANTEE)

↑  
3 MINUTES, 10 SEC INTERVALS

(6)

POLICY: WAIT NO MORE THAN 30 SEC AT A HOME



GETS TO RING ALL 5 HOMES IN REMAINING 3'

$$U^{\pi}(b) = \sum_s \vec{\alpha}_s^T \vec{b}$$

$$\approx [0, \theta]^T [0.15, 0.85] + \begin{matrix} \text{HOUSE \#0} \\ (\text{LIGHT=ON}) \end{matrix}$$

$$[0, \theta]^T [0.15, 0.85] + \begin{matrix} \text{HOUSE \#1} \\ (\text{LIGHT=ON}) \end{matrix}$$

$$[0, \theta]^T [0.95, 0.05] + \begin{matrix} \text{HOUSE \#3} \\ (\text{LIGHT=OFF}) \end{matrix}$$

$$[0, \theta]^T [0.95, 0.05] + \begin{matrix} \text{HOUSE \#4} \\ (\text{LIGHT=OFF}) \end{matrix}$$

$$[0, \theta]^T [0.95, 0.05] + \begin{matrix} \text{HOUSE \#5} \\ (\text{LIGHT=OFF}) \end{matrix}$$

$$x \cdot [0, (1-\theta)^t \cdot \theta]^T [0.15, 0.85] + \begin{matrix} \text{WAITS AT} \\ \text{LIGHT=ON} \end{matrix}$$

$$y \cdot [0, (1-\theta)^t \cdot \theta]^T [0.95, 0.05] + \begin{matrix} \text{WAITS AT} \\ \text{LIGHT=OFF} \end{matrix}$$

~~APPROXIMATE = 0 AS  $\theta \gg (1-\theta)^t \theta$~~

$$U^{\pi}(b) \approx 0.85 \theta \times 3 + 0.05 \cdot \theta \times 2$$

$$U^{\pi}(b) \approx 1.85 \theta = 0.93$$

## ⑦ POMDP MODEL

### 12 STATE VARIABLES:

$[h]$  - IS SOMEONE AT A PARTICULAR HOME: 5 BINARY

$[r]$  - RUNG DOOR AT A PARTICULAR HOME: 5 BINARY

$[l]$  - CURRENT LOCATION

$[t]$  - TIMESTEPS REMAINING

### 5 ACTIONS:

- WAIT: =

- NEXT DOOR:  $\leftarrow \uparrow \rightarrow \downarrow$

### REWARD:

- VISIT SAME HOME AS BEFORE:  $-\infty$

ASSUMPTION: NO

KNOWLEDGE TRANSFER  
ABOUT PREVIOUS CONDO  
VISITS

### FAST INFORMED BOUND:

→ HOW MANY  $\alpha$  VECTORS:

5 (ONE PER ACTION)

- DIMENSIONALITY OF  $\alpha$  VECTORS: # OF STATES

$$\text{STATE} = [h, r, l, t]^T$$

$$|S| = |h| \times |r| \times |l| \times |t|$$

### @ SCENARIO OF Q5.

$$|h| = 5 \times 2 \text{ GIVEN}$$

$$|r| = 5 \times 2 \text{ GIVEN}$$

$$|l| = 5 \text{ (HOMES LEFT)}$$

$$|t| = 18 \text{ (3 MINUTES LEFT)} \\ \text{10 SEC INCREMENTS)}$$

$$|S| = 10 \times 10 \times 5 \times 18$$

$$|S| = 9000$$

16/19

⑧

## ENUMERATED 18-STEP CONDITIONAL PLANS

$\alpha_p$  ALPHA VECTOR (FROM COMPUTING VALUE)  
FOR EACH POSSIBLE  
INITIAL PLAN STATE

RELATIONSHIP?

$$\max_P \vec{\alpha}_P^T b \leq \max_{\alpha \in \Gamma_{FB}} \vec{\alpha}^T b$$

$U^*(b)$

EXACT OPTIMAL  
VALUE FUNCTION

(PSPACE COMPLETE  
COMPLEXITY)

$U^F(b)$

APPROXIMATE FAST  
INFORMED BOUND  
VALUE FUNCTION

$O(|A|^2 |s|^7 |\alpha|)$

⑨

WITH FAST INFORMATION BOUND

$O(|A|^2 |S|^2 |\mathcal{O}|)$  AT EACH ITERATION

TEMPTING AS QMPD WOULD BE

$O(|A|^2 |S|^2)$  AT EACH ITERATION

ANNA WOULD NOT HAVE FULL OBSERVABILITY

AT EVERY NEXT STEP

⇒ SARSOP WOULD BE MY RECOMMENDATION

BECAUSE SHE WOULD ONLY HAVE SOLVE  
FEW  
FOR THE BELIEF STATES REACHABLE BY  
THE OPTIMAL POLICY IN THE 3 MINUTES  
REMAINING.

**From:** Pablo Rodriguez Bertorello <pablo@1now.ai>  
**Sent:** Monday, December 4, 2017 8:56 PM  
**To:** usa5105@fedex.com  
**Subject:** FIRST

#####
##### TRICK OR TREAT? CAT EATS LILY
#####
##### THE HUNGRY CAT NEAR DOOR WAITS
#####
##### WITH TEETH GUARDING HOLY GATES
#####
##### ADORNED IN MASTER'S LOVE SHE SMILES
#####
##### BEEN WAITING TO EAT LILY FOR MILES
#####

10 RHYME

### ### REFERENCES

```
# POMDPs.jl API http://juliapomdp.github.io/POMDPs.jl/latest/api/#POMDPs.state\_index
# Tutorial Tiger POMDP http://nbviewer.jupyter.org/github/sisl/POMDPs.jl/blob/master/examples/Tiger.ipynb
# Explicit POMDP http://juliapomdp.github.io/POMDPs.jl/latest/explicit/
# Tiger problem PDF https://www.cs.rutgers.edu/%7Emlittman/papers/aij98-pomdp.pdf
# Tiger problem PPT https://www.techfak.uni-bielefeld.de/~skopp/Lehre/STdKI\_SS10/POMDP\_tutorial.pdf
# POMDP Simulator https://github.com/JuliaPOMDP/POMDPToolbox.jl
```

### ### POMDPs.jl

```
# The following may be performed prior and separately by own JuliaCommand.jl
import all POMDPs
POMDPs.add("QMDP")          # Design Under Uncertainty 6.4.1
POMDPs.add("SARSOP")
using QMDP
using SARSOP
using POMDPToolbox
using POMDPModels
```

### # TODO: remove

```
# Pkg.clone("https://github.com/sisl/POMDPs.jl.git")
# Pkg.add("POMDPs")
# using POMDPs
# Pkg.add("SARSOP")
# POMDPs.add("QMDP")          # Design Under Uncertainty 6.4.1
# POMDPs.add("POMDPToolbox")   # implements discrete belief updating
# Pkg.add("Distributions")
# import all POMDPs
# using SARSOP
# using POMDPToolbox
```

### ### PROBLEM

```
# In the tiger POMDP, the agent is tasked with escaping from a room. There are two doors leading out of the room.
# Behind one of the doors is a tiger, and behind the other is sweet, sweet freedom.
# If the agent opens the door and finds the tiger, it gets eaten (and receives a reward of -100).
# If the agent opens the other door, it escapes and receives a reward of 10.
# The agent can also listen. Listening gives a noisy measurement of which door the tiger is hiding behind.
# Listening gives the agent the correct location of the tiger 85% of the time.
# The agent receives a reward of -1 for listening.
```

### ### PARTIALLY OBSERVABLE MARKOV DECISION PROBLEM

```
# definition:
# S: State
# A: Action
# Omega: Observation space
# O: Observation function
# T: Transition
# R: Reward
```

PABLO RODRIGUEZ BERTORELLO

19/19