# CDT Summer School

- Introductions (10 min)
  - No email/surfing web/texting
  - Find Github
    - Find your breakout group!

- RL & MRTs (1 hour, 20 min)
  - Lecture 30 min
  - Breakout in groups + Discussion (20 min)
  - Lecture 30 min

1

QR code for https://github.com/StatisticalReinforcementLearningLab/Stat-ML-CDT-2023/tree/main

(Introduce RL and Heartsteps): 3 Hours, 9.30 to 12.30 Monday morning
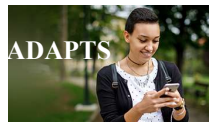
# CDT Summer School

- Break (10 min)

- Deeper dive into bandit algorithms en route to building HeartSteps RL algorithm
  - Lecture 30 min
  - Breakout in groups + Discussion (15 min)
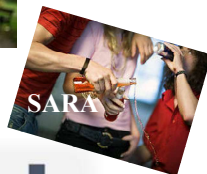  - Lecture & Coding 35 min

2

https://github.com/StatisticalReinforcementLearningLab/Stat-ML-CDT-2023/tree/main/precourse

**Reinforcement Learning & Micro-Randomization**

Susan A Murphy

&

Raaz Dwivedi

some characteristics of digital health in clinical research settings

large no. of stakeholders with differing data needs

need to contribute to behavioral science

low signal to noise ratio with limited data -->intermixing behavioral science with data science

References: (see http://people.seas.harvard.edu/~samurphy/research.html)

Lecture/Practicum/Discussion 1: An introduction to micro-randomized trials, reinforcement learning and bandit algorithms for use in digital intervention development.

References:

1.a Liao,P., Klasnja, P., Tewari, P., Murphy, S.A. Sample Size Calculations for Micro-Randomized Trials in mHealth, Statistics in Medicine. 2016 May 30;35(12):1944-71. [2015 Dec 28 Epub ahead of print] PubMed PMID: 26707831, PMCID PMC4848174.

1.b Russo, D., Van Roy, B., Kazerouni, A., Osband, I., Wen, Z. A Tutorial on Thompson Sampling, https://arxiv.org/abs/1707.02038, Foundations and Trends® in Machine Learning 11 (1), 1-96

1.c Tewari, A. and Murphy, S.A., From Ads to Interventions: Contextual Bandits in Mobile Health, (2017) Mobile Health Sensors, Analytic Methods, and Applications , Springer International Publishing AG 2017, J.M. Rehg et al. (eds.), DOI

10.1007/978-3-319-51394-2_4, pgs 495-518

1.d Hadad Vitor, A Hirshberg David, Zhan Ruohan, Wager Stefan, Athey Susan. Confidence intervals for policy evaluation in adaptive experiments, https://www.pnas.org/doi/10.1073/pnas.2014602118

1.e Kelly Zhang, Lucas Janson, Susan Murphy. Inference for batched bandits. Advances in neural information processing systems 33, 9818-9829, 2020 https://arxiv.org/abs/2002.03217

# To Think About

- What is a Micro-Randomized Trial?

- How are Micro-Randomized Trials and Reinforcement Learning connected?

- What are the two types of policies?!

- What are the two elements of an online RL algorithm?

4

## HeartSteps (PI Klasnja)

Goal: Develop a mobile activity coach for individuals who are at high risk of coronary artery disease

Three iterative studies:

- V1: 42-day micro-randomized study. Study involves sedentary people
- V2/V3: 90-day + 270-day micro-randomized using an online RL algorithm. Study involves people who have Stage 1 Hypertension.

5

N=37 in V1   Data at https://github.com/klasnja/HeartStepsV1
N=42 from V2 and N=49 from V3 users—total =91

blood pressure that falls in the stage 1 hypertension range

The systolic pressure is 140 to 159 mm Hg or your diastolic pressure is 90 to 99 mm Hg

Changing your lifestyle can go a long way toward controlling high blood pressure.

Eating a heart-healthy diet with less salt

Getting regular physical activity

Maintaining a healthy weight or losing weight if you're overweight or obese

Limiting the amount of alcohol you drink
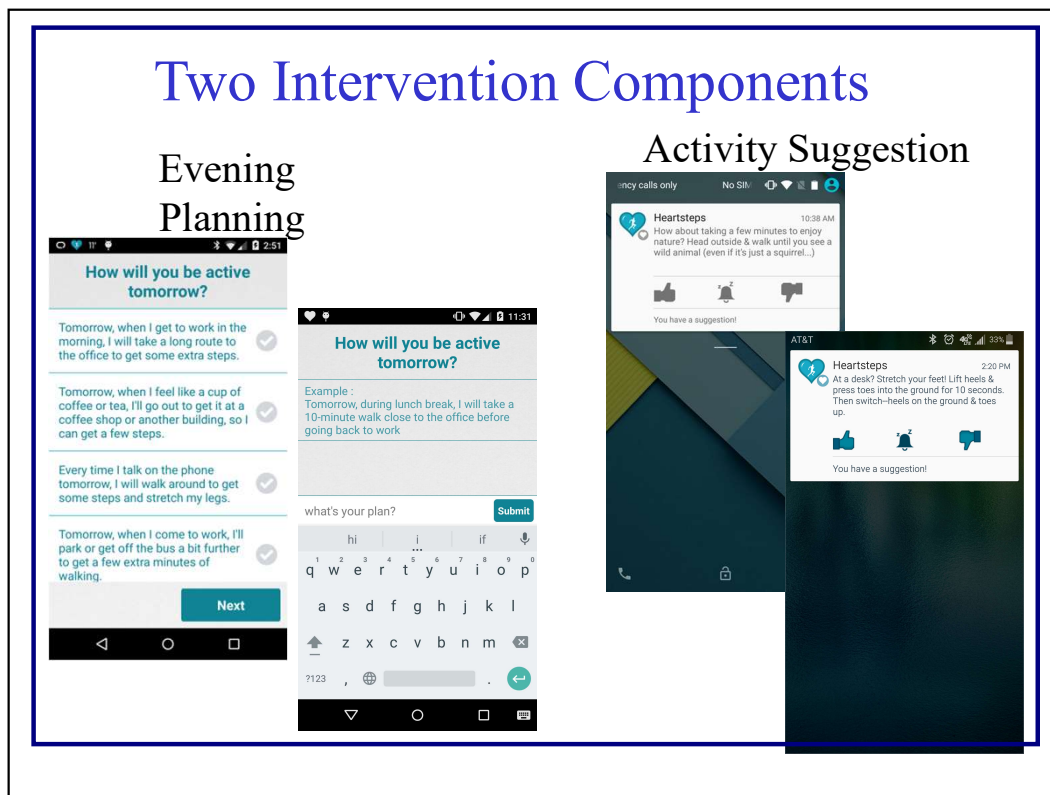
# HeartSteps State

Individual's current state provided via data from:

<u>Wearable band</u> → activity and sleep quality;

<u>Smartphone sensors</u> → busyness of calendar, location, weather, app usage;

<u>Self-report</u> → stress, user burden

6

**Two Intervention Components**

Evening Planning

Activity Suggestion

Activity suggestions are to help in your automated processing

Evening planning is to help in your reflective processing (controlled behavior)—this was in V1; changed to a weekly intervention in V2/V3

There are many intervention components that make up a mobile health intervention. We only experiment with a few.

Example Intervention components
•Whether to provide an intervention or whether to prompt self-monitoring
•How to deliver an intervention option (via a message on wearable, smartphone notification, SMS)
•"Provide nothing" option

Reminders

Suggestions, tips, motivational messages

Prompts to set goals, complete self-report…

Rewards for goal attainment

Recommendations
Reach out recommendation (contact a friend)
Behavioral strategies (exercise;  stay in locations that are supportive of change)
Cognitive strategies (relaxation; reframing)
Motivational messages (reasons for behavior change; barriers for change);

Setting goals; modifying goals
Feedback (often with visualization: fish; flower; garden)
Distractions (game, music, etc.)


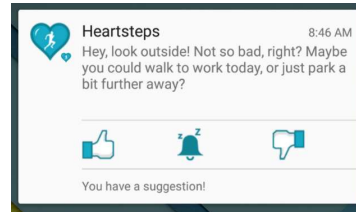Can use sensing and user modeling to determine right delivery time

Don't rely on user's awareness of times of need or remembering to access

But…

High burden

One of many components (good morning message, anti-sedentary message, self-monitoring…)

morning, mid-day, mid-afternoon, early evening, after dinner

Reward: Frequently the actions are primarily designed to have a near-term effect on the individual. E.g. Help then manage current craving/stress, help them manage or be aware of the impact of their social setting on their craving/stress
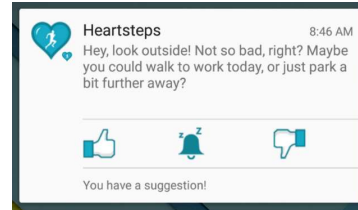
# Micro-randomized Trial

- Micro-randomized trial = each user is randomized many times = sequential experimentation

- Randomization may use Reinforcement Learning (RL)

9

# HeartSteps V1

- 210 decision times per user: 5 time points per day



- Treatment (Action): Contextually tailored activity suggestion (deliver or not deliver); if user is available then randomize with probability 0.6

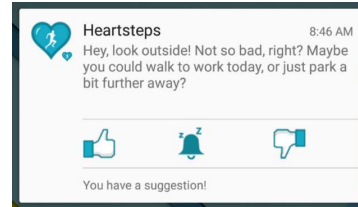- Outcome (Reward): 30-minute step count following each time point.

10

Define availability at a decision time (available if definitely not driving a vehicle, not too active recently, has not turned off delivery in HeartSteps settings…)

"micro-randomized" actions in V1
Suggestion tailored to location, weather, day of week, time of day

# HeartSteps V2/V3

- 450 or > 1000 decision times per user: 5 time points per day



- Treatment (Action): Contextually tailored activity suggestion (deliver or not deliver); randomization via an online RL algorithm

- Outcome (Reward): 30-minute step count following each time point.

11

Suggestion tailored to location, weather, day of week, time of day

# Data for Online Reinforcement Learning (RL) Algorithm

- User's data at time $t : \{S_t, A_t, R_{t+1}\}$

- $t$: decision time (5 times per day)
- $S_t$: State accrued up to/including time $t$
- $A_t$: Action at $t^{\text{th}}$ time ($a = 1$ deliver, $a = 0$ do not deliver) if state permits.
- $R_{t+1}$: Reward accrued after time $t$ and up to time $t + 1$ (log of the 30 min step count)

12

# Online RL Algorithm

- Data at time $t : \{S_t, A_t, R_{t+1}\}$

- For $t = 1$ to $T$ do:
  - Algorithm receives state $S_t$
  - Algorithm selects action $A_t$
  - Algorithm receives reward $R_{t+1}$
- End For

13

This is often called a "contextual bandit environment" ($S_t$ does not depend on prior actions and the reward, $R_{t+1}$ only depends on time t state, action)

We only get to observe one of the two rewards (one of two labels) depending on which action is selected.

This is what makes RL different from

Supervised ML in which you usually see all labels

Unsupervised ML in which you see no labels

# Simple Decision-Making Setting

- Policy $\pi: \mathcal{S} \to \mathcal{P}_{\{0,1\}}$ (distribution on $\{0,1\}$)

- Value of policy $\pi$:

$$V(\pi) = \mathbb{E}\left[\sum_{a\in\{0,1\}} \pi(a|S_t)\, R_{t+1}(a)\right]$$

15

Policy $\pi$ may be degenerate (AKA, select a particular action with probability =1 for each state)

$\mathbb{E} == E_{\{S_t, R_{t+1}(1), R_{t+1}(0)\}\sim\mathcal{P}}$

In this simple bandit environment setting people often call the value the expected reward.

# Simple Decision-Making Setting

- $V(\pi) = \mathbb{E}\left[\sum_{a \in \{0,1\}} \pi(a|S_t) R_{t+1}(a)\right]$
  - Also written as $V(\pi) = \mathbb{E}_\pi[R_{t+1}]$

- The value depends on the reward function:
$$r(s,a) = \mathbb{E}[R_{t+1}(a)|S_t = s]$$

16

Policy $\pi$ may be degenerate (AKA, select a particular action with probability =1 for each state)

$\mathbb{E} == E_{\{S_t, R_{t+1}(1), R_{t+1}(0)\} \sim \mathcal{P}}$

## Simple Decision-Making Setting

- Reward function:

$$r(s, a) = \mathbb{E}[R_{t+1}(a)|S_t = s \,]$$

- The value:
  - $V(\pi) = \mathbb{E}_\pi[R_{t+1}]$

  $$= \mathbb{E}\left[\sum_{a \in \{0,1\}} \pi(a|S_t)\, R_{t+1}(a)\right]$$

  $$= \mathbb{E}\left[\sum_{a \in \{0,1\}} \pi(a|S_t) r(S_t, a)\right]$$

17

$$V(\pi) = \mathbb{E}_\pi[R_{t+1}]$$
$$= \mathbb{E}\left[\sum_{a \in \{0,1\}} \pi(a|S_t)\, R_{t+1}(a)\right]$$
$$= \mathbb{E}\left[\sum_{a \in \{0,1\}} \pi(a|S_t) r(S_t, a)\right]$$
$$= \mathbb{E}_\pi[r(S_t, A_t)]$$

In this simple bandit environment setting people often call the value the expected reward.

Recall from slide 13 the algorithm only sees $S_t$ (does not see $R_{t+1}(1)$, $R_{t+1}(0)$, prior to selecting the action $A_t$) thus conditional on $S_t$, $A_t$ is independent of $R_{t+1}(1)$, $R_{t+1}(0)$. We obtain:
$$r(s, a) = \mathbb{E}[R_{t+1}(a)|S_t = s \,]$$
$$= \mathbb{E}[R_{t+1}(a)||S_t = s, A_t = a]$$
$$= \mathbb{E}[R_{t+1}|S_t = s, A_t = a]$$

# Simple Decision-Making Setting

- Value of policy $\pi$:

$$V(\pi) = \mathbb{E}\left[\sum_{a \in \{0,1\}} \pi(a|S_t)r(S_t, a)\right]$$

- Evaluation of policy $\pi$?
  - "Highest Value in a comparison class" minus $V(\pi)$

  - Overall optimal policy: $\pi^*(a|s) = 1_{a=\underset{a'}{\mathrm{argmax}}\, r(s,a')}$ [18]

The classical comparison class contains $\pi^*$

However in real life settings this may not be the case due to constraints that have to be imposed on the set of possible policies.

# Back to Online RL Algorithm

- Data at time $t$ : $\{S_t, A_t, R_{t+1}\}$
- History at time $t$:  $H_{t-1} = \{S_j, A_j, R_{j+1}\}_{j=1}^{t-1}$

- For $t = 1$ to $T$ do:
  - Algorithm receives state $S_t$
  - Algorithm selects action $A_t$
  - Algorithm receives reward $R_{t+1}$
- End For

19

# Online RL Algorithm

An online RL algorithm is a sequence of policies,
$$\pi^{\mathcal{L}} = \left\{\pi_t^{\mathcal{L}}\right\}_{t \geq 1}$$

- Each $\pi_t^{\mathcal{L}}$ is a probability distribution on $\{0,1\}$ of form $\pi_t^{\mathcal{L}}(\cdot \mid S_t, H_{t-1})$

- $\pi_t^{\mathcal{L}}$ is used, online, to select $A_t$
  - the distribution of $A_t$ conditional on $S_t, H_{t-1}$ is $\pi_t^{\mathcal{L}}$ [20]

Use of $\pi_t^{\mathcal{L}}$ to select $A_t$ is micro randomization…

History: $H_{t-1} = \left\{S_j, A_j, R_{j+1}\right\}_{j=1}^{t-1}$

# Online RL Algorithm

An online RL algorithm is a sequence of policies,
$\pi^{\mathcal{L}} = \left\{\pi_t^{\mathcal{L}}\right\}_{t \geq 1}$

- Evaluation of $\pi^{\mathcal{L}}$?
  - Regret:

    "best expected cumulative reward in a comparison class" minus $\mathbb{E}_{\pi^{\mathcal{L}}}\left[\sum_{t=1}^{T} R_{t+1}\right]$ for each $T$

21

Raaz will discuss a very nice online decision making alg called a "Linear Thompson-Sampling" Algorithm.

For a great tutorial see Sections 1-6 in

Russo, D., Van Roy, B., Kazerouni, A., Osband, I., Wen, Z. A Tutorial on Thompson Sampling, https://arxiv.org/abs/1707.02038, Foundations and Trends® in Machine Learning 11 (1), 1-96

# Online RL Algorithm

A "regression" perspective to constructing an online decision-making algorithm, $\pi^{\mathcal{L}} = \left\{\pi_t^{\mathcal{L}}\right\}_{t \geq 1}$

- Two Elements
  - <u>Learning Algorithm</u>: An algorithm that estimates parameters in a model for (parts of) $\left\{S_j, A_j, R_{j+1}\right\}_{j=1}^{t-1}$
  - <u>Optimization Algorithm:</u> Uses output of learning algorithm to construct $\pi_t^{\mathcal{L}}$ and select $A_t$

23

In our class we focus on the "regression" perspective. Feels as if we are doing supervised learning even though that is not true…..

The other common perspective is a classification perspective (e.g. policy search)
Many methods combine the two perspectives. Actor-Critic algorithms do this.

# *Example*
## Simple Contextual Bandit

<u>Learning Algorithm</u>: An algorithm that estimates parameters in a model for (parts of) $\{S_j, A_j, R_{j+1}\}_{j=1}^{t}$

- – Algorithm incrementally estimates parameters in a model for $r(s, a) = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$

$R_{t+1}$ is the target

It is via the following that we turn an RL problem into a regression problem (supervised learning):::::::

Recall from slide 13,19 the algorithm only sees $S_t$ (does not see $R_{t+1}(1)$, $R_{t+1}(0)$, prior to selecting the action $A_t$) thus conditional on $S_t$ , $A_t$ is independent of $R_{t+1}(1)$, $R_{t+1}(0)$.  We obtain:

$r(s, a) = \mathbb{E}[R_{t+1}(a) | S_t = s\ ]$
$\qquad = \mathbb{E}[R_{t+1}(a) || S_t = s, A_t = a]$
$\qquad = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$

## *Example*
## Simple Contextual Bandit

Optimization Algorithm: Uses output of learning algorithm to construct $\pi_t^{\mathcal{L}}$ and select $A_t$

- Uses estimators, $\hat{r}(s,1), \hat{r}(s,0),$ and some measure of confidence (estimated standard error, length of estimated confidence interval, posterior variance) to construct $\pi_t^{\mathcal{L}}$
- Uses current state $S_t$ and $\pi_t^{\mathcal{L}}$ to select $A_t$

25

# *Example*
# Simple Contextual Bandit

- Learning Algorithm: An algorithm that estimates parameters in a model for (parts of) $\{S_j, A_j, R_{j+1}\}_{j=1}^{t}$

  - Algorithm incrementally estimates parameters in a model for $r(s,a) = \mathbb{E}[R_{t+1}|S_t = s, A_t = a]$

- Optimization Algorithm: Uses output of learning algorithm to construct $\pi_t^{\mathcal{L}}$ and select $A_t$

  - Uses estimators, $\hat{r}(s,1), \hat{r}(s,0),$ and some measure of confidence (estimated standard error, length of estimated confidence interval, posterior variance) to construct $\pi_{t}^{\mathcal{L}}$
    26

20 min for Discussion!

# Break & Discussion

- What is the difference between a decision-making algorithm policy $\pi^{\mathcal{L}}$ and a policy $\pi$ in the simple decision-making setting?

- How might policy $\pi^{\mathcal{L}}$ be related to a policy $\pi$ in the simple decision-making setting?

- Think about HeartSteps; why might you assume or not assume the simple Bandit environment setting for HeartSteps?

28

# To Think About

- What are Markov Decision Processes?

- What is invariant in the likelihood for an MDP?

- Bellman's Equation & Regression Targets

- What are some implications of the Policy Improvement Theorem?

29

# Markov Decision Process (MDP)

Simplified setting: binary actions, discrete states

Potential states
$$\{S_0, S_1(0), S_1(1), S_2(0), S_2(1), \dots, S_{t+1}(0), S_{t+1}(1) \dots \}$$
drawn from a distribution $\mathcal{P}$

Simplified setting: reward, $R_{t+1}(a) = h(S_{t+1}(a))$ where $h$ is bounded, known

30

# Markovian

Potential states:

$$\{S_0, S_1(0), S_1(1), S_2(0), S_2(1), \dots, S_{t+1}(0), S_{t+1}(1) \dots \}$$

Markovian:

For b ∈ {0,1},

$\{S_{t+1}(0), S_{t+1}(1)\}$ are independent of the past,

$\{S_0, S_1(0), S_1(1), S_2(0), S_2(1), \dots, S_{t-1}(0), S_{t-1}(1)\}$,

conditional on the present state: $S_t(b)$

31

Think carefully –what DAG would represent the potential states with the Markovian property?

# Markovian and Stationary

Potential states:
$$\{S_0, S_1(0), S_1(1), S_2(0), S_2(1), \ldots, S_{t+1}(0), S_{t+1}(1) \ldots\}$$

Markovian and Stationary:

For $a, b \in \{0,1\}, t \geq 0,$

$p(s'|s, a) = P(S_{t+1}(a) = s'|S_0, S_1(0), \ldots, S_{t-1}(1), S_t(b) = s)$

$\Rightarrow p(s'|s, a) = P(S_{t+1}(a) = s'| S_t(b) = s)$

(No dependence on $b$.......; No dependence on $t$....) 32

# MDP

Potential states
$$\{S_0, S_1(0), S_1(1), S_2(0), S_2(1), \ldots, S_{t+1}(0), S_{t+1}(1) \ldots\}$$

- The data collector uses a policy to sample the potential states.
- This policy determines the distribution of $A_0, A_1, \ldots A_t$

33

Markovian, Stationary policy

# Example: Stationary Policy

Potential states
$$\{S_0, S_1(0), S_1(1), S_2(0), S_2(1), \ldots, S_{t+1}(0), S_{t+1}(1) \ldots\}$$

Data collector uses a stationary policy , $\pi: \mathcal{S} \to \mathcal{P}_{\{0,1\}}$ to collect data:

1. Draws $A_0 \sim \pi(\cdot \mid S_0)$ thus $P(A_0 = a \mid S_0 = s) = \pi(a \mid s)$
2. Observes $S_1(A_0)$; define $S_1 \triangleq S_1(A_0)$
3. Draws $A_1 \sim \pi(\cdot \mid S_1)$ thus $P(A_1 = a \mid S_1 = s) = \pi(a \mid s)$
4. Observes $S_2(A_1)$; define $S_2 \triangleq S_2(A_1)$
5. and so on....

34

# MDP

Compete Data:
$$\{S_0, S_1(0), S_1(1), S_2(0), S_2(1), \dots, S_{t+1}(0), S_{t+1}(1), \dots\}$$

Observed Data:
$$\{S_0, A_0, S_1, A_1, S_2, \dots, A_{t+1}, S_{t+1}, \dots\}$$

Recall:

- $R_{t+1} \triangleq h\big(S_{t+1}(A_t)\big) = h(S_{t+1})$
- $r(s, a) \triangleq \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$
  $$= \mathbb{E}[h(S_{t+1}) | S_t = s, A_t = a]$$

35

# MDP

Stationary, Markovian policy, $\pi$.

- Likelihood for observed data, $\{S_j, A_j\}_{j \geq 0}$

  - $p_0(s_0) \prod_{j \geq 0} \left( \pi(a_j|s_j) p(s_{j+1}|a_j, s_j) \right)$
  - $p_0(s_0) = P(S_0 = s_0)$;
  - $p(s'|a, s) = P\left(S_{j+1} = s'|A_j = a, S_j = s\right)$

- Stationary transitions $\implies$ for any $t$,
  $P_\pi\left(S_{j+t} = s'|A_j = a, S_j = s\right)$ is the same for all $j \geq 0$

What do we learn by considering this likelihood?

# Likelihood

- Likelihood for $\{S_t, A_t\}_{t \geq 0}$
  - $p_0(s_0) \prod_{t \geq 0} (\pi(a_t|s_t) p(s_{t+1}|a_t, s_t))$

- Stationary transitions: $p(s'|a, s)$ does not depend on time $t$

"Groundhog Day"

Each time the process visits state $s$, the process starts over.

$$\mathbb{E}_\pi[h(S_{j+t+1})|S_j = s] = \mathbb{E}_\pi[h(S_{j'+t+1})|S_{j'} = s]$$

**Groundhog Day** (1993): Bill Murray stars as a weatherman who gets stuck in a time loop in Punxsutawney, Pennsylvania. He has to relive the same day over and over again until he learns to appreciate life.

Recall

$R_{t+1} \triangleq h(S_{t+1}(A_t)) = h(S_{t+1})$

$(r(s, a) = \mathbb{E}[R_{t+1}|S_t = s, A_t = a])$

# MDP

$$R_{t+1} \triangleq h\big(S_{t+1}(A_t)\big) = h(S_{t+1})$$

$$\mathbb{E}_\pi\big[h(S_{j+t+1})|S_j = s\big] = \mathbb{E}_\pi\big[h(S_{j'+t+1})|S_{j'} = s\big]$$
$$\Longrightarrow$$
$$\mathbb{E}_\pi\big[R_{j+t+1}|S_j = s\big] = \mathbb{E}_\pi\big[R_{j'+t+1}|S_{j'} = s\big]$$

for all $t, j \geq 0$

38

# MDP

Consider stationary, Markovian policies:
- Discounted sum of rewards
  - $\sum_{t=0}^{\infty} \gamma^t R_{t+1}$  (discounted horizon, $0 < \gamma < 1$)

- Value of policy $\pi$

$$V^{\pi}(s) = \mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t R_{t+1} | S_0 = s]$$
$$= \mathbb{E}_{\pi}\left[\sum_{t=0}^{\infty} \gamma^t R_{j+t+1} | S_j = s\right] \text{ (for all } j \geq 0)$$

39

There are other types of sums of rewards one can consider…

See Puterman, Markov Decision Processes, 2005  Ch. 5, 6 for a proof of why we can restrict attention to stationary, markovian policies.

Note that
$V^{\pi}(s) = \mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t R_{k+t+1} | S_k = s]$ for all k  "groundhog day"

# Bellman's Equation

- Stationary, Markovian policy $\pi$

$$V^\pi(s)$$
$$= \mathbb{E}_\pi[\sum_{t=0}^\infty \gamma^t R_{t+1} \,|S_0 = s]$$
$$= \mathbb{E}_\pi[R_1 + \gamma \sum_{t=1}^\infty \gamma^{t-1} R_{t+1} \,|S_0 = s]$$
$$= \mathbb{E}_\pi[R_1 + \gamma \mathbb{E}_\pi[\sum_{t=1}^\infty \gamma^{t-1} R_{t+1} \,|S_1]|S_0 = s]$$
$$= \mathbb{E}_\pi[R_1 + \gamma \mathbb{E}_\pi[\sum_{t=0}^\infty \gamma^t R_{1+t+1} \,|S_1]|S_0 = s]$$
$$= \mathbb{E}_\pi[R_1 + \gamma V^\pi(S_1)|S_0 = s]$$

40

$V^\pi(s) = \mathbb{E}_\pi[\sum_{t=0}^\infty \gamma^t R_{k+t+1} \,|S_k = s]$ for all k  "groundhog day"

# Bellman's Equation

Bellman's equation for policy $\pi$

$$V^\pi(s) = \mathbb{E}_\pi[R_1 + \gamma V^\pi(S_1)|S_0 = s]$$
$$\quad\quad = \mathbb{E}_\pi[R_{t+1} + \gamma V^\pi(S_{t+1})|S_t = s]$$

$V^\pi(s)$
$$= \sum_a \pi(a|s)\{\mathbb{E}[R_{t+1} + \gamma V^\pi(S_{t+1})|S_t = s, A_t = a]\}$$

Q-function
- $Q^\pi(s, a) \triangleq \mathbb{E}[R_{t+1} + \gamma V^\pi(S_{t+1})|S_t = s, A_t = a]$

41

Q-function – Quality of action a in state s

Also called the Action-Value function

Q-function – Quality of rewards accrued by first taking action a in state s and then follow policy pi thereafter.

$$Q^\pi(s, a) = \mathbb{E}_\pi[\sum_{t=0}^{\infty} \gamma^t R_{t+1} |S_0 = s, A_t = a]$$

# Bellman's Equation

Bellman's equation for policy $\pi$

$$V^\pi(s)$$
$$= \sum_a \pi(a|s)\mathbb{E}[R_{t+1} + \gamma V^\pi(S_{t+1})|S_t = s, A_t = a]$$

$$Q^\pi(s, a) \triangleq \mathbb{E}[R_{t+1} + \gamma V^\pi(S_{t+1})|S_t = s, A_t = a]$$

$$V^\pi(s) = \sum_a \pi(a|s)Q^\pi(s, a)$$

$$Q^\pi(s, a) = \mathbb{E}_\pi[R_{t+1} + \gamma Q^\pi(S_{t+1}, A_{t+1})|S_t = s, A_t = a]$$

Q-function – Quality of rewards accrued by first taking action a in state s and then follow policy pi thereafter.

$$
\begin{aligned}
V^\pi(s) &= \sum_a \pi(a|s)\mathbb{E}[R_{t+1} + \gamma V^\pi(S_{t+1})|S_t = s, A_t = a] \\
&= r(s, a) + \gamma \sum_a \pi(a|s)\mathbb{E}[V^\pi(S_{t+1})|S_t = s, A_t = a] \\
&= r(s, a) + \gamma \mathbb{E}_\pi[V^\pi(S_{t+1})|S_t = s, A_t = a]
\end{aligned}
$$

# Bellman's Equation

- Bellman's equation for the Value Function

$$V^\pi(s) = \mathbb{E}_\pi[R_{t+1} + \gamma V^\pi(S_{t+1})|S_t = s]$$

- Bellman's equation for the Action Value Function (Q-Function)

$$Q^\pi(s,a) = \mathbb{E}_\pi[R_{t+1} + \gamma Q^\pi(S_{t+1}, A_{t+1})|S_t = s, A_t = a]$$

$Q^\pi(s,a)$ is the expected discounted rewards accrued by initially taking action $a$ and then following $\pi$ thereafter.

43

Q-function – Quality of rewards accrued by first taking action a in state s and then follow policy pi thereafter.

$$V^\pi(s) = \sum_a \pi(a|s)\mathbb{E}[R_{t+1} + \gamma V^\pi(S_{t+1})|S_t = s, A_t = a]$$
$$= r(s,a) + \gamma \sum_a \pi(a|s)\mathbb{E}[V^\pi(S_{t+1})|S_t = s, A_t = a]$$
$$= r(s,a) + \gamma \mathbb{E}_\pi[V^\pi(S_{t+1})|S_t = s, A_t = a]$$

$$Q^\pi(s,a) = \mathbb{E}_\pi[R_{t+1} + \gamma Q^\pi(S_{t+1}, A_{t+1})|S_t = s, A_t = a]$$
$$= r(s,a) + \gamma \mathbb{E}[\sum_{a'} \pi(a'|s)Q^\pi(S_{t+1}, a')|S_t = s, A_t = a]$$

## Regression Targets

$$V^\pi(s) = \mathbb{E}_\pi[R_{t+1} + \gamma V^\pi(S_{t+1})|S_t = s]$$
$$\implies$$

Time $t$ Value target: $R_{t+1} + \gamma V^\pi(S_{t+1})$

$$Q^\pi(s,a) = \mathbb{E}_\pi[R_{t+1} + \gamma Q^\pi(S_{t+1}, A_{t+1})|S_t = s, A_t = a]$$
$$\implies$$

- Time $t$ Q target: $R_{t+1} + \gamma Q^\pi(S_{t+1}, A_{t+1})$

44

Regression refers to conditional expectations…

These are common targets if you want to learn the value of a particular policy…

Another time t (less noisy) Q target:
$R_{t+1} + \gamma \sum_{a'} \pi(a'|s) Q^\pi(S_{t+1}, a')$
Why less noisy?

# MDP

- Value of policy $\pi$:

$$V^\pi(s) = \mathbb{E}_\pi\left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \mid S_0 = s\right]$$

- Evaluation of policy $\pi$?
  - "Best Value in a comparison class" minus $V^\pi$

  - Common choice is overall Value maximizer:

$$\pi^* = \operatorname*{argmax}_{\pi \in \Pi} V^\pi(s)$$

45

$\Pi$ is the set of markovian, stationary policies of which there is only $2^{|\mathcal{S}|}$. That is, if there are $|\mathcal{S}|$ discrete states and 2 actions, then the number of possible policies is $2^{|\mathcal{S}|}$ ; for each state there are two possible actions.

Why does the optimal policy not depend on initial state s?

Why does there exist a markovian, stationary optimal policy in the set of all policies?

See Puterman, Markov Decision Processes, 2005 Ch. 5 (do not need to consider history dependent policy, only consider markovian policies), Ch 6 for a proof of why we can restrict attention to stationary policies and the optimal policy does not depend on initial state.

# Optimal Policy, Optimal Q-Function

Common choice is overall Value maximizer:
$$\pi^* = \underset{\pi}{\text{argmax}} \, V^\pi(s)$$

"Optimal" policy, $\pi^*$ is $\pi^*(a|s) = 1_{a=\underset{a'}{\text{argmax}} \, Q^{\pi^*}(s,a')}$

Equivalently,
$$\pi^*(s) = \underset{a}{\text{argmax}} \, Q^{\pi^*}(s,a) \text{ where } \pi^*(s) \in \{0,1\}$$

46

A common target if you are aiming to learn the optimal policy in an unconstrained stetting.

Optimal policy is degenerate….

Why is $\pi^*(s) = \underset{a}{\text{argmax}} \, Q^{\pi^*}(s,a)$?

# Bellman's Optimality Equation

"Optimal" policy, $\pi^*$ is $\pi^*(a|s) = 1_{a=\underset{a\prime}{\text{argmax}}\, Q^{\pi^*}(s,a\prime)}$ ????

Bellman's equation:
$$Q^{\pi^*}(s,a) =$$
$$\mathbb{E}\left[R_{t+1} + \gamma \sum_{a\prime} \pi^*(a'|S_{t+1})Q^{\pi^*}(S_{t+1}, a') \,|S_t = s, A_t = a\right]$$
$$\Longrightarrow$$
Bellman's Optimality Equation
$$Q^{\pi^*}(s,a) = \mathbb{E}\left[R_{t+1} + \gamma \max_{a\prime} Q^{\pi^*}(S_{t+1}, a') \,|S_t = s, A_t = a\right]$$

47

Why is $\pi^*(s) = \underset{a}{\text{argmax}}\, Q^{\pi^*}(s,a)$?

We know this due to the policy improvement theorem on the following slides. Policy improvement theorem also implies that there is an optimal policy that is deterministic....

## Optimal Regression Target

"Optimal" policy, $\pi^*$ is $\pi^*(a|s) = 1_{a=\underset{a'}{\text{argmax }} Q^{\pi^*}(s,a')}$

$$Q^{\pi^*}(s,a) = \mathbb{E}\left[R_{t+1} + \gamma \max_{a'} Q^{\pi^*}(S_{t+1}, a') \,|\, S_t = s, A_t = a\right]$$

Time $t$ Q target: $R_{t+1} + \gamma \max_{a'} Q^{\pi^*}(S_{t+1}, a')$

48

Why is $\pi^*(s) = \underset{a}{\text{argmax }} Q^{\pi^*}(s,a)$?

We know this due to the policy improvement theorem on the following slides. Policy improvement theorem also implies that there is an optimal policy that is deterministic….

# Bellman's Optimality Equation

"Optimal" policy, $\pi^*$ is $\pi^*(a|s) = 1_{a=\text{argmax}\, Q^{\pi^*}(s,a')}$

Optimality equations are:

$$Q^{\pi^*}(s,a) = \mathbb{E}\left[R_{t+1} + \gamma \max_{a'} Q^{\pi^*}(S_{t+1}, a') \,|\, S_t = s, A_t = a\right]$$

and

$$V^{\pi^*}(s) = \sum_{a'} \pi^*(a'|s) Q^{\pi^*}(s, a')$$
$$= \max_{a'} Q^{\pi^*}(s, a')$$
$$= \max_{a'} \mathbb{E}[R_{t+1} + \gamma V^{\pi^*}(S_{t+1}) | S_t = s, A_t = a']$$

49

Recall

$Q^{\pi}(s,a) \triangleq \mathbb{E}[R_{t+1} + \gamma V^{\pi}(S_{t+1}) | S_t = s, A_t = a]$ for all stationary, markovian policies $\pi$

Why is $\pi^*(s) = \text{argmax}_a Q^{\pi^*}(s,a)$?

We know this due to the policy improvement theorem on the following slides. Policy improvement theorem also implies that there is an optimal policy that is deterministic....

## Policy Improvement Theorem

Suppose $\pi$ and $\pi'$ are two stationary, Markovian policies for which

$\sum_a \pi'(a|s)Q^\pi(s,a) \geq \sum_a \pi(a|s)Q^\pi(s,a)$ for all $s$,

then

$V^{\pi'}(s) = \sum_a \pi'(a|s)Q^{\pi'}(s,a) \geq \sum_a \pi(a|s)Q^\pi(s,a) = V^\pi(s)$

for all $s$

Common choice is $\pi'(a|s) = 1_{a=\underset{a'}{\operatorname{argmax}} Q^\pi(s,a')}$

50

To see how this implies that $\pi^*(a|s) = 1_{a=\operatorname{argmax}_{a'} Q^{\pi*}(s,a')}$, let $\pi=\pi^*$ and consider $\pi'(a|s) = 1_{a=\underset{a'}{\operatorname{argmax}} Q^{\pi*}(s,a')}$ aka, put $\pi'(s) = \underset{a}{\operatorname{argmax}} Q^{\pi*}(s,a)$,

then the policy improvement theorem implies that $V^{\pi'}(s) \geq V^{\pi^*}(s)$ so $\pi'$ is an optimal policy.

Chapter 4 in Sutton and Barto

# Policy Improvement Theorem

Suppose $\pi$ and $\pi'$ are two stationary, Markovian policies for which
$\sum_a \pi'(a|s)Q^\pi(s,a) \geq \sum_a \pi(a|s)Q^\pi(s,a)$ for all $s$,
then $V^{\pi'}(s) \geq V^\pi(s)$ for all $s$

$$Q^\pi(s,a) = \mathbb{E}_\pi[R_{t+1} + \gamma Q^\pi(S_{t+1}, A_{t+1})|S_t = s, A_t = a]$$

Time $t$ Q target for policy improvement:
$R_{t+1} + \gamma Q^\pi(S_{t+1}, A_{t+1})$ where $A_{t+1} \sim \pi$

51

Another time t (less noisy) Q target:
$R_{t+1} + \gamma \sum_{a'} \pi(a'|s)Q^\pi(S_{t+1}, a')$
Why less noisy?

# Policy Improvement Theorem Proof

Suppose $\sum_a \pi'(a|s)Q^\pi(s,a) \geq \sum_a \pi(a|s)Q^\pi(s,a)$ for all $s$,

Step 1. Show that for any $t$,
$$Q^\pi(S_t, A_t) \leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma Q^\pi(S_{t+1}, A_{t+1})|S_t, A_t]$$

(see notes to this slide!)

52

$$
\begin{aligned}
Q^\pi(S_t, A_t) &= \mathbb{E}_\pi[R_{t+1} + \gamma Q^\pi(S_{t+1}, A_{t+1})|S_t, A_t] \\
&= \mathbb{E}[R_{t+1} + \gamma \mathbb{E}_\pi[Q^\pi(S_{t+1}, A_{t+1})|S_{t+1}]|S_t, A_t] \\
&= \mathbb{E}[R_{t+1} + \gamma \sum_a \pi(a|S_{t+1})Q^\pi(S_{t+1}, a)|S_t, A_t] \\
&\leq \mathbb{E}[R_{t+1} + \gamma \sum_a \pi'(a|S_{t+1})Q^\pi(S_{t+1}, a)|S_t, A_t] \\
&= \mathbb{E}[R_{t+1} + \gamma \mathbb{E}_{\pi'}[Q^\pi(S_{t+1}, A_{t+1})|S_{t+1}]|S_t, A_t] \\
&= \mathbb{E}_{\pi'}[R_{t+1} + \gamma Q^\pi(S_{t+1}, A_{t+1})|S_t, A_t]
\end{aligned}
$$

# Policy Improvement Theorem Proof

For all $t$, prior slide shows that
$$Q^\pi(S_t, A_t) \leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma Q^\pi(S_{t+1}, A_{t+1})|S_t, A_t]$$

2. Show that
$$Q^\pi(S_0, A_0) \leq \mathbb{E}_{\pi'}\left[\sum_{j=0}^{t} \gamma^j R_{j+1} + \gamma^{t+1} Q^\pi(S_t, A_t)|S_0, A_0\right]$$
implies
$$Q^\pi(S_0, A_0) \leq \mathbb{E}_{\pi'}\left[\sum_{j=0}^{t+1} \gamma^j R_{j+1} + \gamma^{t+2} Q^\pi(S_{t+1}, A_{t+1})|S_0, A_0\right]$$

53

To obtain 2. just use the result from prior slide: $Q^\pi(S_t, A_t) \leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma Q^\pi(S_{t+1}, A_{t+1})|S_t, A_t]$ inside the expectation on the right-hand side.

## Policy Improvement Theorem Proof

3. Then by induction

$$Q^\pi(S_0, A_0) \leq \mathbb{E}_{\pi'}\left[\sum_{j=0}^{t} \gamma^j R_{j+1} + \gamma^{t+1} Q^\pi(S_t, A_t)|S_0, A_0\right]$$

holds for all $t$.

Let $t \to \infty$,

$$Q^\pi(S_0, A_0) \leq \mathbb{E}_{\pi'}\left[\sum_{j=0}^{\infty} \gamma^j R_{j+1}|S_0, A_0\right]$$
$$= Q^{\pi'}(S_0, A_0)$$

54

$R_{j+1}$ is a.s. bounded and $\gamma$ is in (0,1). Thus the limit $t \to \infty$ can be taken inside of the expectation.

Policy improvement theorem also states that if,

$\sum_a \pi'(a|s)Q^\pi(s,a) \geq \sum_a \pi(a|s)Q^\pi(s,a)$ for all $s$, with a strict inequality for one state

Then

$V^\pi(s) \leq V^{\pi'}(s)$ for all states, with a strict inequality for at least one state.

# Online RL Algorithm

An online decision-making algorithm is a sequence of policies, $\pi^{\mathcal{L}} = \left\{\pi_j^{\mathcal{L}}\right\}_{j \geq 1}$

- Many RL algorithms use the policy improvement theorem in some way to construct the $\pi_j^{\mathcal{L}}$'s

55

# RL Conceptual Idea

1. Estimate $Q^{\pi}$ for the current policy, $\pi$
   1. $Q^{\pi}(s, a) = \mathbb{E}_{\pi}[R_{t+1} + \gamma Q^{\pi}(S_{t+1}, A_{t+1})|S_t = s, A_t = a]$

2. Find a $\pi'$ for which
   $$\sum_a \pi'(a|s)Q^{\pi}(s, a) \geq \sum_a \pi(a|s)Q^{\pi}(s, a) \text{ for all } s$$
   1. One choice is $\pi'(a|s) = 1_{a = \underset{a'}{\operatorname{argmax}} Q^{\pi}(s, a')}$

3. Conclude that $\pi'$ is a better policy than $\pi$,
   $$V^{\pi'}(s) = \sum_a \pi'(a|s)Q^{\pi'}(s, a) \geq \sum_a \pi(a|s)Q^{\pi}(s, a) = V^{\pi}(s)$$

56

Time $t$ Q target for policy improvement: $R_{t+1} + \gamma Q^{\pi}(S_{t+1}, A_{t+1})$ where $A_{t+1} \sim \pi$

# Conceptual Idea

1. Estimate $Q^\pi$ for the current policy, $\pi$ $\qquad \left(\pi \leftarrow \pi_j^{\mathcal{L}}\right)$

2. Find a $\pi'$ for which

   $$\sum_a \pi'(a|s) Q^\pi(s,a) \geq \sum_a \pi(a|s) Q^\pi(s,a) \text{ for all } s$$

3. Conclude that $\pi'$ is a better policy than $\pi$,
   $$V^{\pi'}(s) \geq V^\pi(s)$$

4. Use $\pi'$ to select actions. $\qquad \left(\pi_{j+1}^{\mathcal{L}} \leftarrow \pi'\right)$

57

# To Think About!

- What is the difference between a decision-making algorithm policy $\pi^{\mathcal{L}}$ and a policy $\pi$ in the MDP setting?

- How might policy $\pi^{\mathcal{L}}$ be related to a policy $\pi$ in the MDP setting?

- How do the targets differ between the bandit environment and the Markovian environment?

58