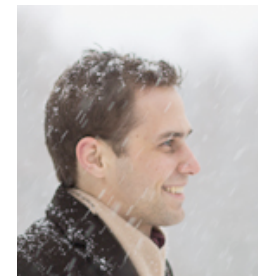# Learning the Preferences of Ignorant, Inconsistent Agents

**Owain Evans (Oxford)**, Andreas Stuhlmueller (Stanford),
Noah Goodman (Stanford)

# 1. Motivation for learning human preferences

- Scientific (economics, psychology): how do people value work vs. leisure, short-term vs. long-term, country vs. friends & family?

- Machine learning (applications): recommendation (movie, job, dating), create tailored content.

- Machine learning (long-term goal): the more systems **understand** our preferences, the more they can help us make **high stakes** decisions in **novel** circumstances.

**Lisa Larter**
Edit Profile

**FAVORITES**
- 📋 News Feed
- 💬 Messages — 49
- 7 Events — 20+
- 🖼 Photos
- 🔍 Browse
- 📇 Ads Manager
- 📰 Social Fixer News
- 📄 Done For You Pages — 20+
- 🔴 Lisa Larter
- 🌐 The Pilot Project 2... — 1
- 🚩 The Pilot Project G... — 2

**PAGES**
- ⛰ Exclusive Associat... — 1
- 🖼 Tanner the Little ... — 4
- 🌳 Branching Out
- eWomenNetwork ... — 6
- eWomenNetwork ... — 20+
- eWomenNetwork Orang...
- eWomenNetwork ... — 6
- 📣 Pages Feed — 20+
- 👍 Like Pages — 20+

**GROUPS**
- 🍏 TPP ATP
- 🌐 WIBWS
- 📋 20VIC Marketing — 4
- 🟢 Create Group...

**FRIENDS**
- ⭐ Close Friends — 20+

**APPS**
- 🎮 Games — 8
- 🎮 Games Feed — 20+

**INTERESTS**
- 📋 Clients
- 📗 Add Interests...

---

📝 **Update Status**   📷 **Add Photos/Video**

What's on your mind?

SORT ▾

**Natalie Deschamps** shared VR-Zone's photo.

Skateboard baby stroller that comes with brakes and handlebars for steering

👍 Like VR-Zone for more amazing stuff

Like · Comment · Share · 4 minutes ago · 🌐

👍 Angela Azaria likes this.

Angela Azaria Hilarious.
2 minutes ago · Like

Write a comment... 📷

**Christine Tripp** at The Centurion Conference & Event Center
Looking forward to seeing Cara!
Like · Comment · Share · 11 minutes ago in Ottawa · 🌐

Write a comment... 📷

⌄

**Sara Karissa**

---

🎂 **Julie Azizan** and 3 others
7 **2 events** today

**Trending**                    Learn More

↗ **York University:** 2 women injured after shooting at York University

↗ **International Women's Day:** International Women's Day –– How Empowering One Impacts Many

↗ **Wayne Gretzky:** MacKinnon breaks Gretzky record in Avalanche win over Red Wings

▾ See More

# 2. Learning preferences with IRL

**Inverse Reinforcement Learning (AI) / Structural Estimation (Econ):**

- Unsupervised learning, assumed model is MDP, POMDP, RL.

- Learn from sequences of choices in complex environments (cf. Netflix)

- Learn utility/reward function not policy: enduring cause not contingent effects.

- People act on their preferences without ability to report them quantitatively (driving skill, detailed vacation plan)

# 3. The problem of systematic error

- IRL: infer preferences from observed actions ... assuming human fits (MDP/POMDP) model up to random (softmax) errors.

- But human make **systematic** errors! Person smokes every day but regrets it.

- Behavioral economics (hyperbolic discounting, Prospect Theory)

- Bounded cognition (forgetting, limited computational ability, etc.)

# 4. Learning from ignorant, inconsistent agents

Our approach:

1. build flexible generative models to capture a range of biases and cognitive bounds (while maintaining tractability)

2. jointly infer **biases** (or lack thereof) and **preferences** from behavior

3. if successful, can help humans overcome biases

# 5. Human bias: Time inconsistency

- Intuition: tonight you want to rise early but tomorrow you want to sleep in.

- Most prominent bias: addiction, procrastination, impulsiveness, will-power / pre-commitment.

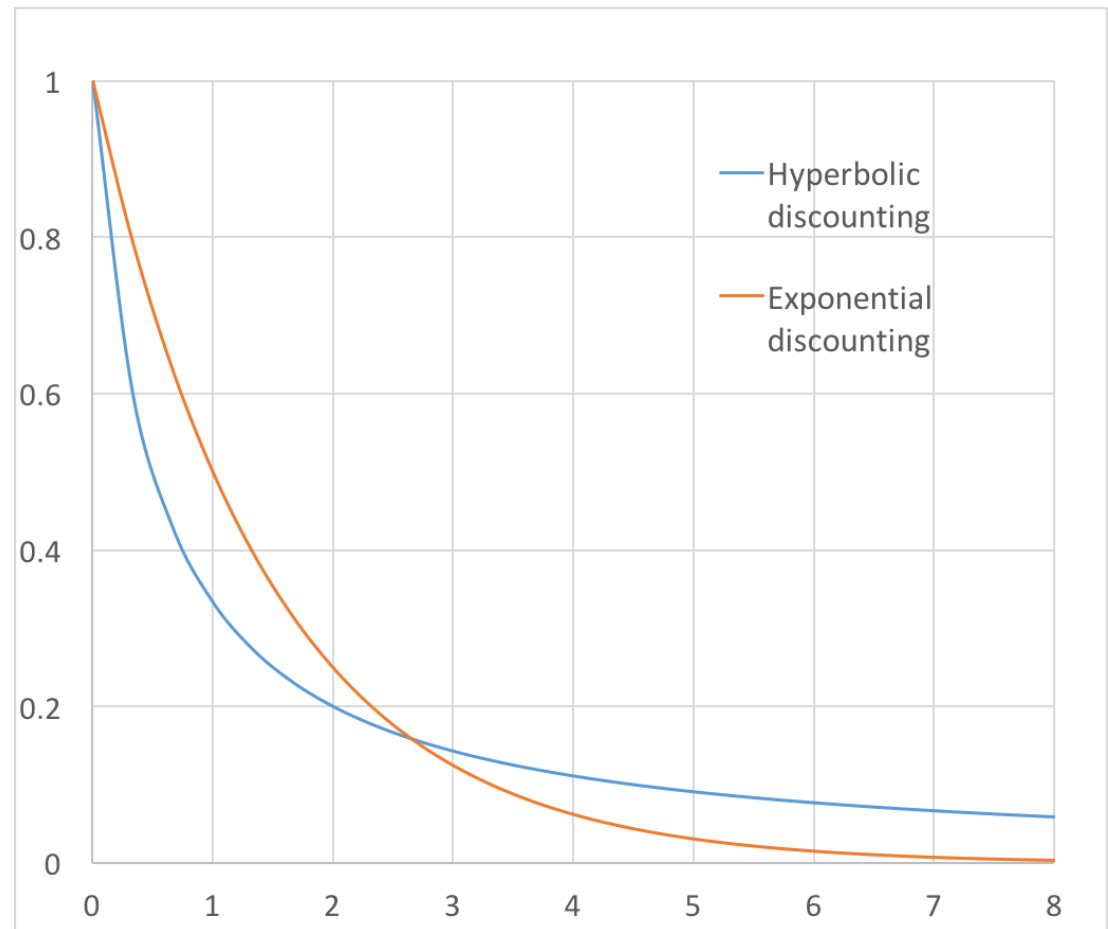- Formally, any non-exponential discounting implies time-inconsistency.
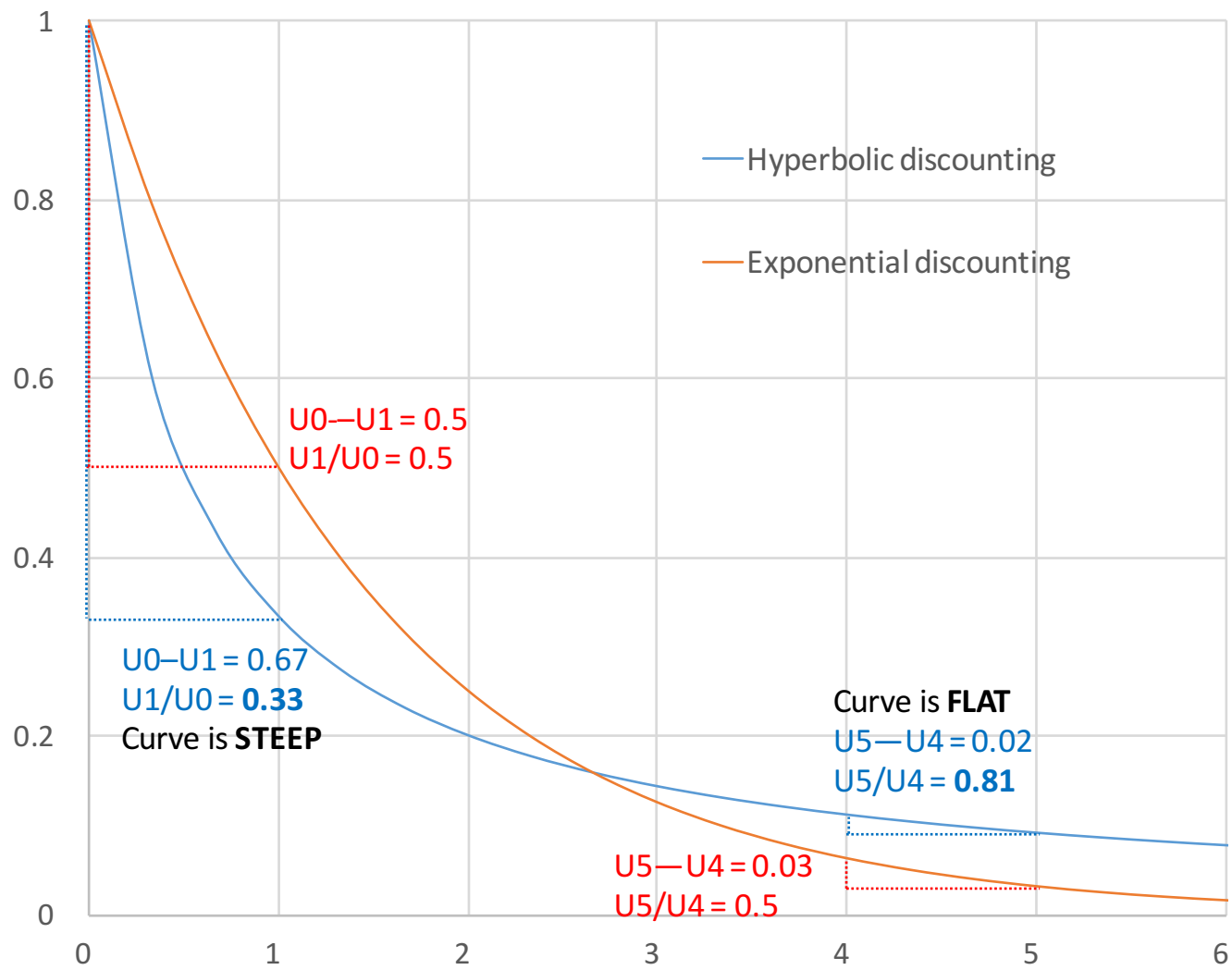
# 5. Human bias: Time inconsistency

**Hyperbolic discounting**

Discount factor = $1/(1+kt)$

At t=0, you prefer $80 at t=8 to $70 at t=7 (curve **shallow**)

At t=7, you re-evaluate and prefer $70 now to $80 tomorrow (curve **steep)**.

# 5. Model for biased agent

**MDP model:**
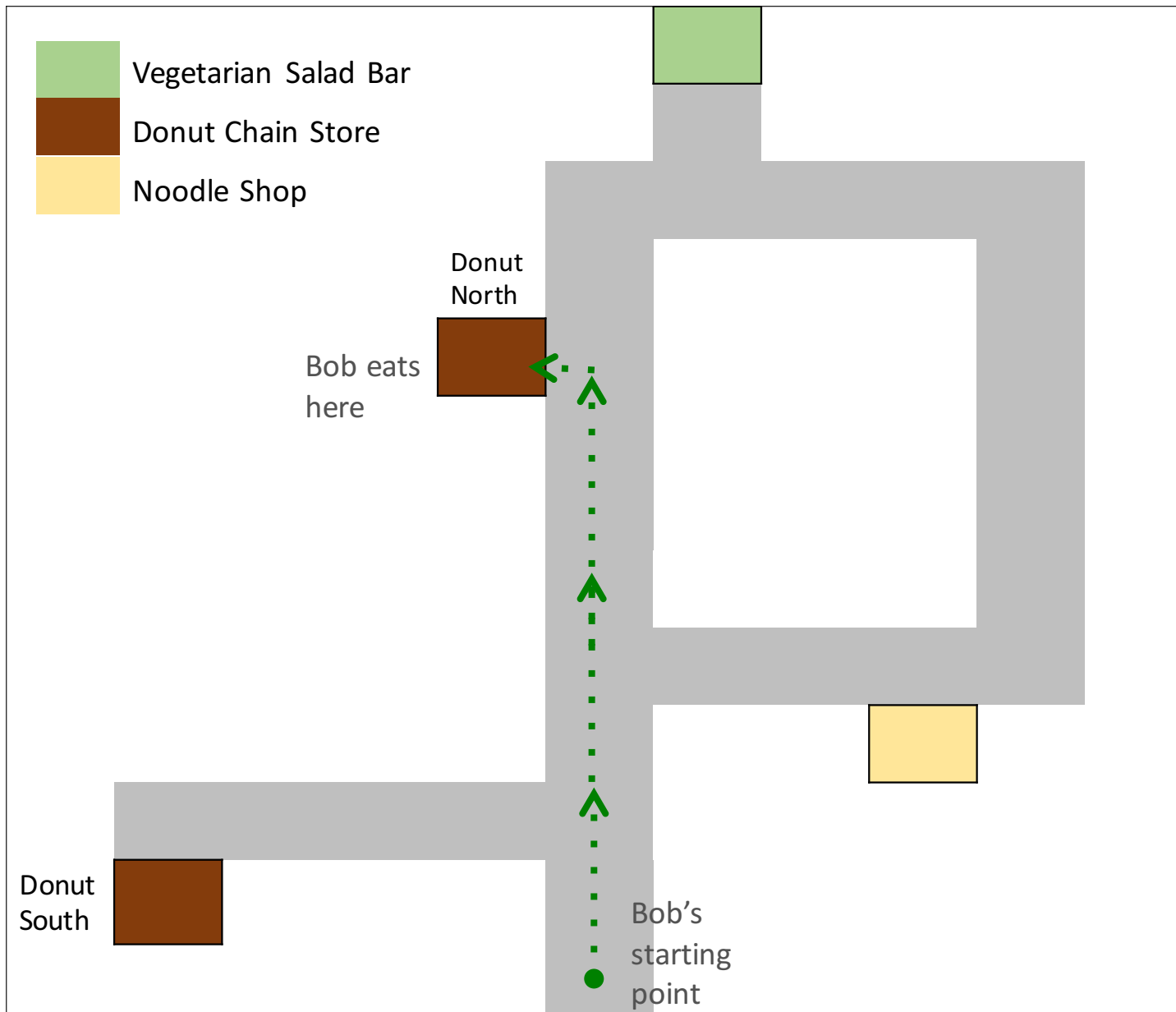$$\text{EU}_s[a] = U(s, a) + \mathop{\mathbb{E}}_{s', a'}[\text{EU}_{s'}[a']]$$

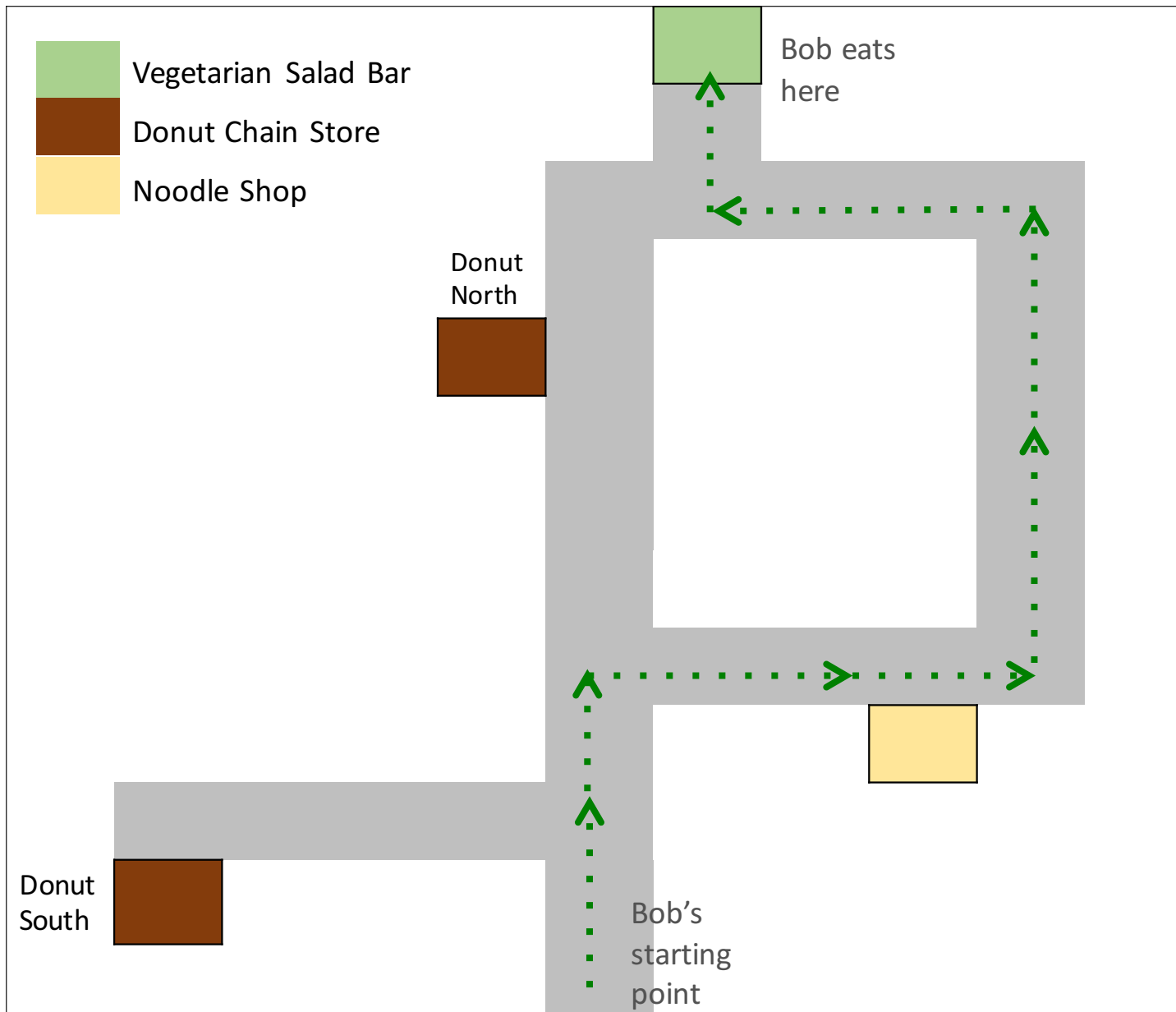$$\text{with } s' \sim T(s, a) \text{ and } a' \sim C(s')$$

**MDP + Hyberbolic discounting** (variable *d* for "delay" measures how far in the future the action *a* would take place):

$$\text{EU}_{s,d}[a] = \frac{1}{1 + kd} U(s, a) + \mathop{\mathbb{E}}_{s', a'}[\text{EU}_{s', d+1}[a']]$$

# 6. Goal for examples and experiments

- Show that ignoring biases (assuming optimality) leads to mistakes in learning preferences

- Mistakes occur in simple, uncontrived, everyday scenarios.

Vegetarian Salad Bar

Donut Chain Store

Noodle Shop

Donut North

Bob eats here

Donut South

Bob's starting point

Vegetarian Salad Bar

Donut Chain Store

Noodle Shop

Bob eats here

Donut North

Donut South

Bob's starting point

# 5. Model for biased agent - NAIVE

**MDP model:**
$$\mathrm{EU}_s\left[a\right] = U(s,a) + \mathop{\mathbb{E}}_{s',a'}\left[\mathrm{EU}_{s'}\left[a'\right]\right]$$
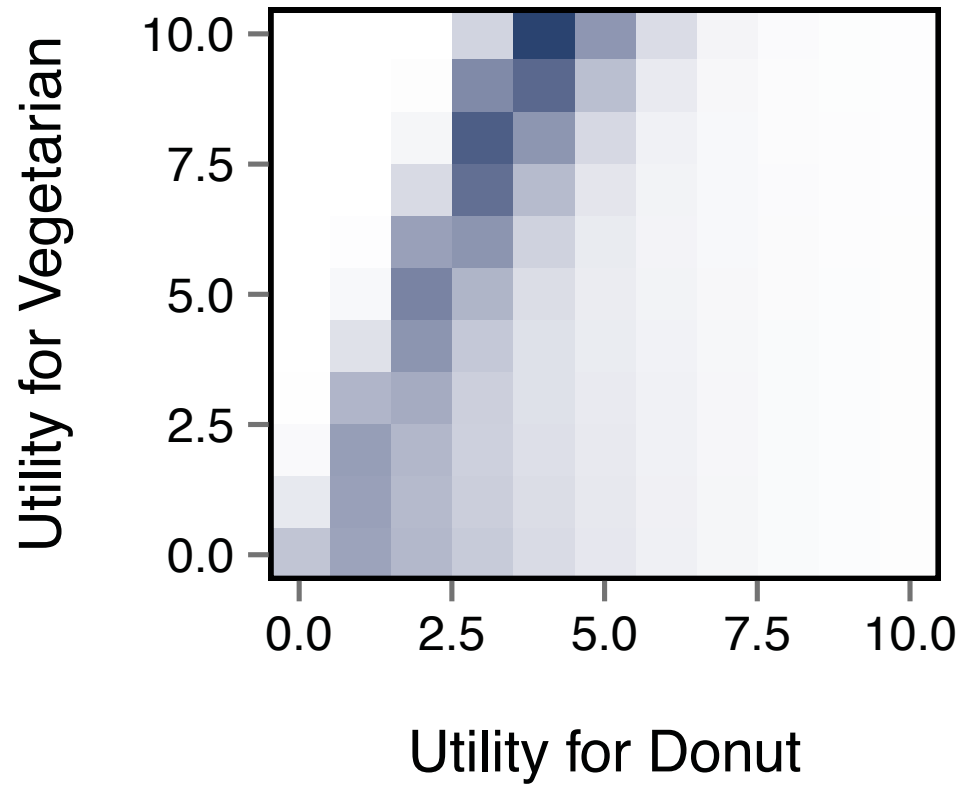
with $s' \sim T(s,a)$ and $a' \sim C(s')$

**MDP + Hyberbolic discounting** (variable *d* for "delay" measures how far in the future the action *a* would take place):

$$\mathrm{EU}_{s,d}\left[a\right] = \frac{1}{1+kd}U(s,a) + \mathop{\mathbb{E}}_{s',a'}\left[\mathrm{EU}_{s',d+1}\left[a'\right]\right]$$

$$a' \sim C(s', d+1)$$

# 5. Model for biased agent - SOPHISTICATED

**MDP model:**
$$\mathrm{EU}_s\left[a\right] = U(s,a) + \mathop{\mathbb{E}}_{s',a'}\left[\mathrm{EU}_{s'}\left[a'\right]\right]$$

with $s' \sim T(s,a)$ and $a' \sim C(s')$

**MDP + Hyberbolic discounting** (variable *d* for "delay" measures how far in the future the action *a* would take place):
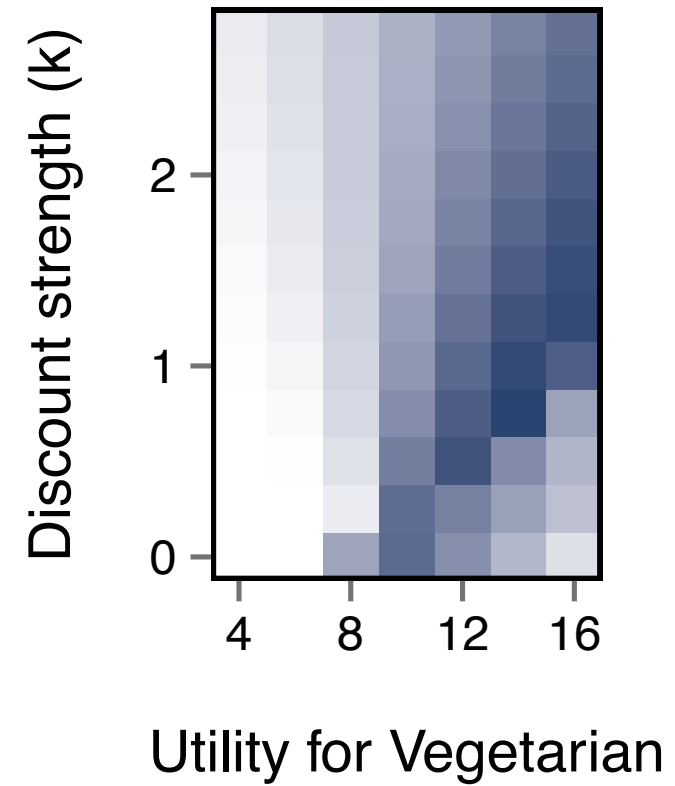
$$\mathrm{EU}_{s,d}\left[a\right] = \frac{1}{1+kd}U(s,a) + \mathop{\mathbb{E}}_{s',a'}\left[\mathrm{EU}_{s',d+1}\left[a'\right]\right]$$

$$a' \sim C(s',0)$$

Naive

Sophisticated

# 6. Model for biases agent: Procrastination



do nothing $u = 0$

promise $u = -\epsilon$

do work $u = -1$

help friend $u = +R$
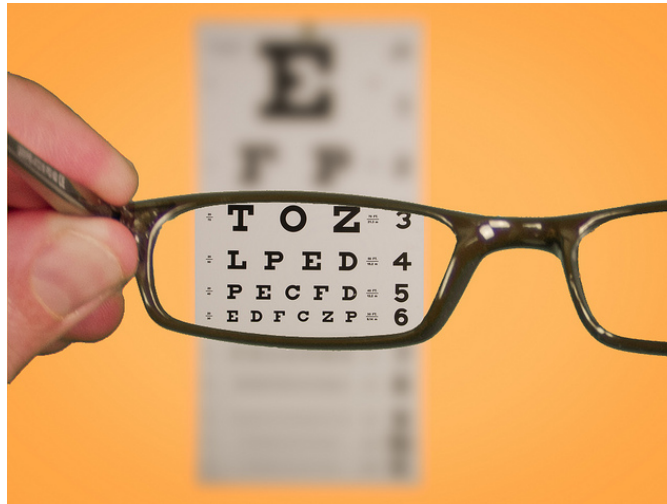
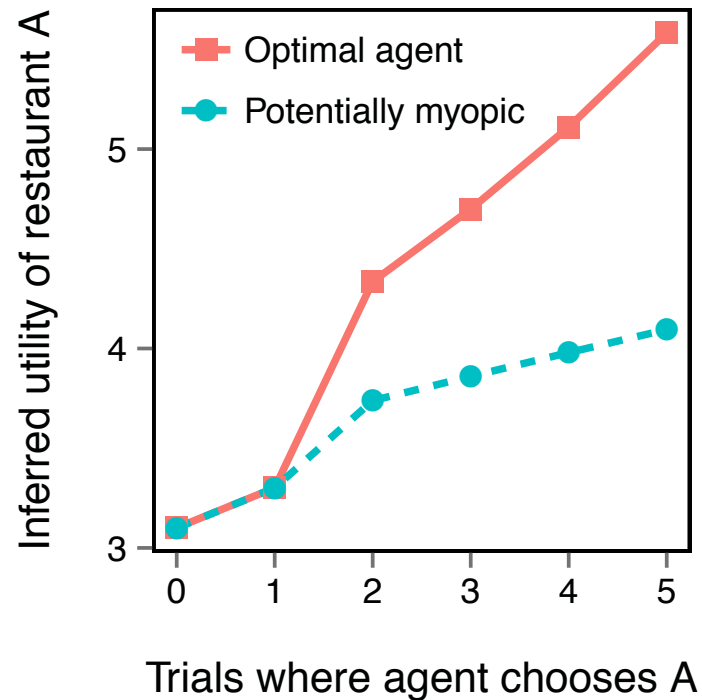# 6. Model for biased agent: Procrastination

# 7. Model for biased agent: Myopia

- **Simple myopia (near sighted)**: ignore any rewards or costs after time *k1 > 0* (even though you'll still be alive).

- **Bounded Value-of-Information:** ignore the value of information gained after time *k2 > 0* (even though you will still get benefits from information).
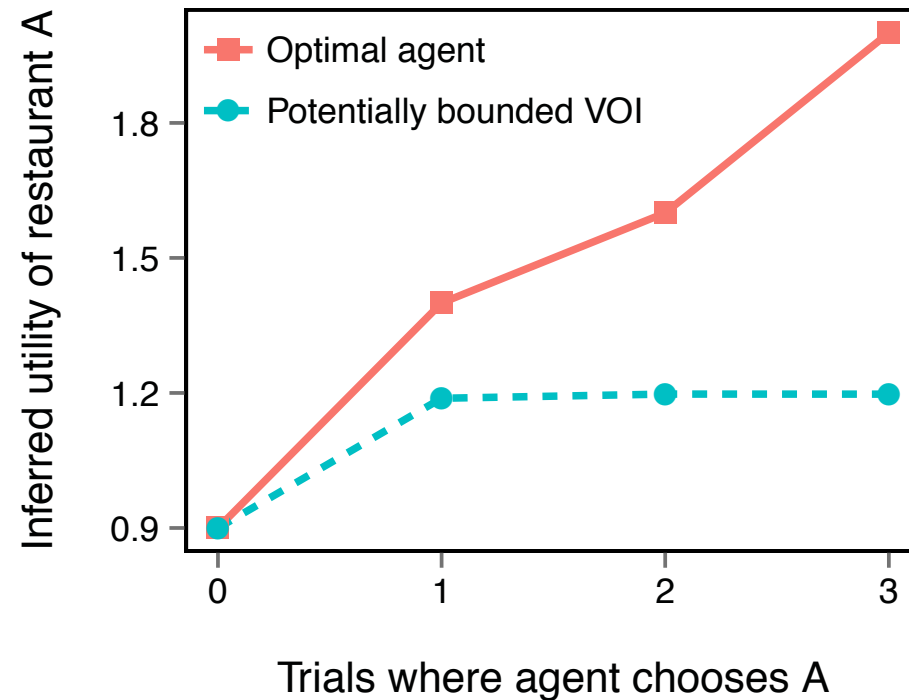
# 7. Model for biased agent: Myopia



Myopic planning

Bounded VOI

# agentmodels.org

Interactive, online tutorial and open-source library for constructing this kind of model (Work in progress).

Main sections:

- Agent models for one-player sequential problems (MDPs, POMDPs, RL), where agent can be biased

- Inference (IRL) for a large space of possible agents

- Multi-agent interactions: coordination, group preferences.

Tom's *decision rule* is to take action $a$ that maximizes utility, i.e., the action

$$\arg\max_{a \in A} U(T(s, a))$$

In WebPPL, we can implement this utility-maximizing agent as a function `maxAgent` that takes a state $s \in S$ as input and returns an action. For Tom's choice between restaurants, we assume that the agent starts off in a state `"default"`, denoting whatever Tom does before going off to eat. The program directly translates the decision rule above using the higher-order function `argMax`.

```
// Choose to eat at the Italian or French restaurants
var actions = ['italian', 'french'];

var transition = function(state, action){
  return (action === 'italian') ? 'pizza' : 'steak frites';
};

var utility = function(state){
  return (state === 'pizza') ? 1 : 0;
};

var maxAgent = function(state){
  return argMax(
    function(action){
      return utility(transition(state, action));
    },
    actions);
};

print("Agent chooses: " + maxAgent("default"));
```
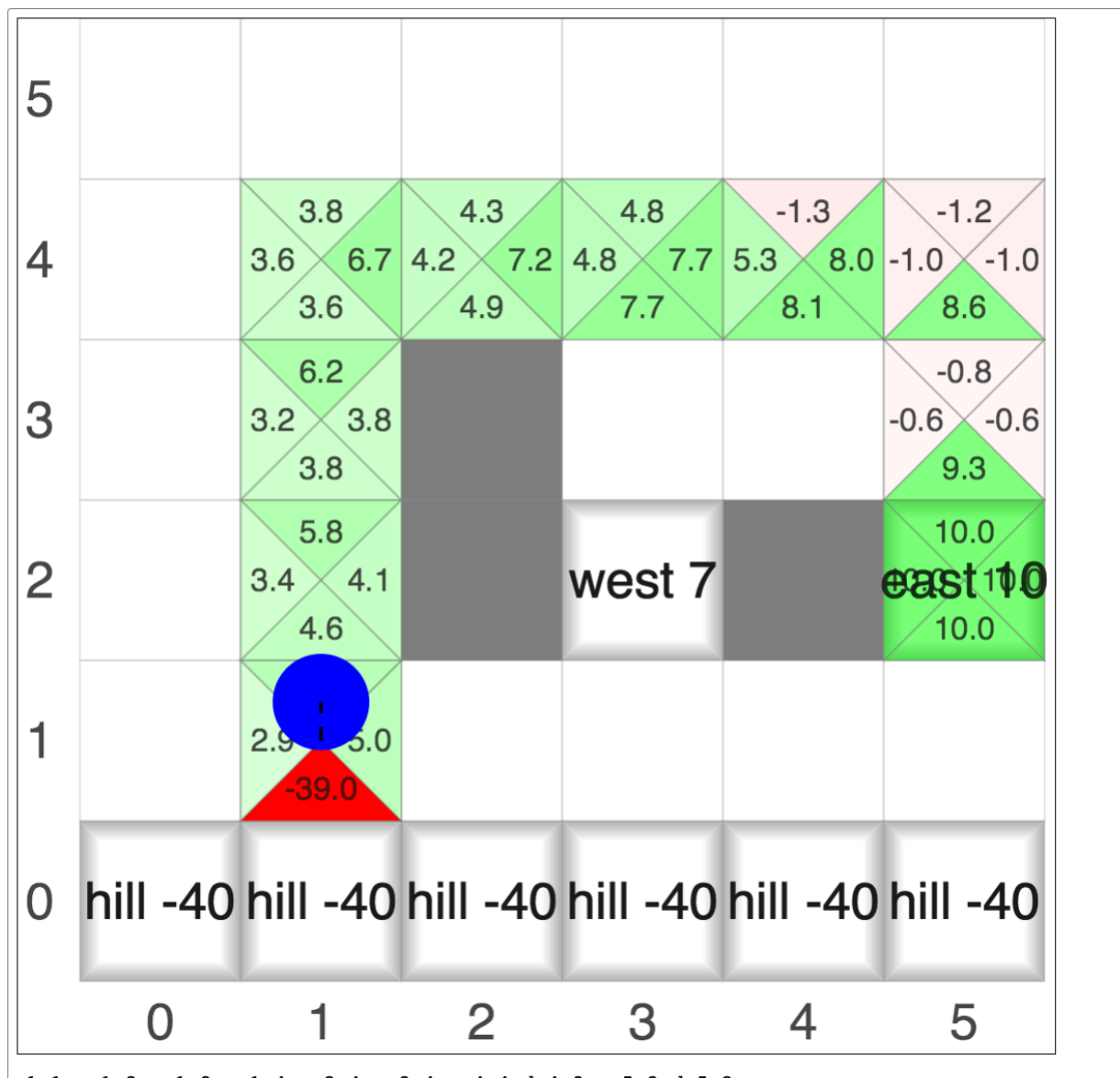
> run

```
Agent chooses: french
```

# Acknowledgments