

Regularized Approaches for Soft MDP

Pu Yang
School of Mathematical Sciences, Peking University

2022.11.23

Table of contents

- 1 Introduction
- 2 Soft MDP
 - Soft Q-learning
 - Soft Actor-Critic
- 3 Learn a Sparse Policy
- 4 Application: Mobile Health

- 1 Introduction
- 2 Soft MDP
- 3 Learn a Sparse Policy
- 4 Application: Mobile Health

Markov Decision Processes (MDP) $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle^1$

- \mathcal{S} - a set of states; $s \in \mathcal{S}$ - a state
- \mathcal{A} - a set of actions; $a \in \mathcal{A}$ - an action (we only consider discrete case)
- P - transition probability function: $\mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$
- R - reward function: $\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$
- γ - discounting factor for future rewards

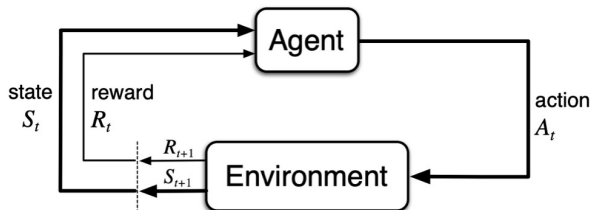


Figure 1: The agent–environment interaction in a Markov decision process.

¹Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

Bellman Equation

- policy probability function: $\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$
- Value function: $V^\pi(s) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^k r_t \mid s_0 = s \right]$
- Q function: $Q^\pi(s, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^k r_t \mid s_0 = s, a_0 = a \right]$

Bellman equation

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_{a \sim \pi} [Q(s_t, a) \mid s_t = s] \\ &= \mathbb{E}_\pi [r_t + \gamma V^\pi(s_{t+1}) \mid s_t = s] \end{aligned}$$

$$\begin{aligned} Q^\pi(s, a) &= \mathbb{E} [r_t + \gamma V^\pi(s_{t+1}) \mid s_t = s, a_t = a] \\ &= \mathbb{E} [r_t + \gamma \mathbb{E}_{a \sim \pi} Q^\pi(s_{t+1}, a) \mid s_t = s, a_t = a] \end{aligned}$$

Bellman Optimality Equation

We aim to find the optimal policies by solving

$$\pi^* = \arg \max_{\pi} V^{\pi}(s)$$

- The optimal value function $V^* = V^{\pi^*}$ is unique solution of

$$V^*(s) = \mathcal{T}V^*(s)$$

where \mathcal{T} is called Bellman Operator defined as

$$\mathcal{T}V(s) = \max_{\pi} \mathbb{E}_{\pi} [r_{t+1} + \gamma V(s_{t+1}) \mid s_t = s]$$

- The optimal Q function $Q^* = Q^{\pi^*}$ satisfies

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')$$

- The optimal policies (denoted as π^*) may not be unique

$$\pi^*(a \mid s) = 0, \quad \text{if } Q(s, a) < \max_a Q(s, a)$$

but π^* often is a deterministic policy which puts all probability mass on one action

Deterministic Policy vs. Stochastic Policy

In some cases, we might actually prefer to learn stochastic policies.
For example

- Multi-goal environment (see Fig. 2)
- Exploration (like DFS vs. BFS) (see Fig. 3)
- Robustness in the face of unstationary environment (may be future work)

Multi-goal Environment

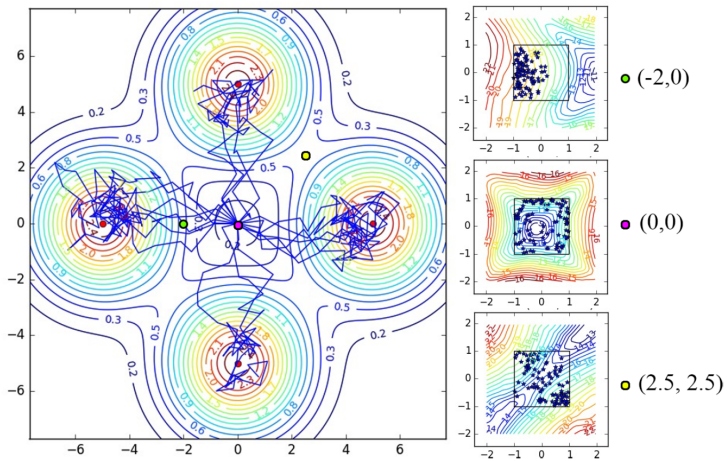


Figure 2: Illustration of the 2D multi-goal environment.

Exploration

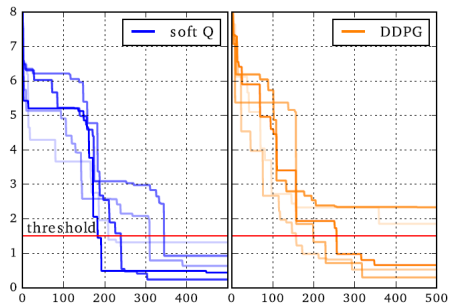
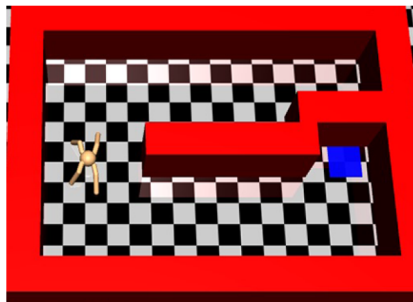


Figure 3: Comparison of soft Q-learning and DDPG the quadrupedal robot maze task.

- 1 Introduction
- 2 Soft MDP**
- 3 Learn a Sparse Policy
- 4 Application: Mobile Health

Preliminary 1: Maximum Entropy Reinforcement Learning

- standard reinforcement learning objective

$$\pi_{\text{std}}^* = \arg \max_{\pi} \sum_t \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}} [r(s_t, a_t)]$$

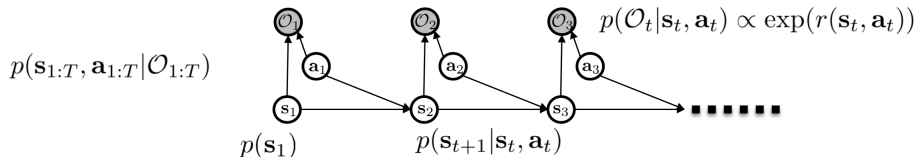
- augments the reward with an entropy regularization term

$$\pi_{\text{MaxEnt}}^* = \arg \max_{\pi} \sum_t \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}} [r(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot | s_t))],$$

where α is the temperature parameter that controls the stochasticity of the policy, and the entropy is $\mathcal{H}(P) = \mathbb{E}_{x \sim P}[-\log P(x)]$

A Derivation from Probabilistic Graphical Model

We could consider control problems as the following temporal probabilistic graphical model (PGM)



where s_1, s_2, \dots and a_1, a_2, \dots are hidden variables, representing states and actions, respectively; O_1, O_2, \dots are observed binary variables, which indicate whether the corresponding state and action are optimal.

A Derivation from Probabilistic Graphical Model

- Probabilistic inference problem: figure out the trajectory distribution given an optimal policy, i.e. $p(s_{1:T}, a_{1:T} \mid \mathcal{O}_{1:T})$
- We assume $O_{1:T} = \mathbf{1}$, which indicates the probability of a trajectory given that it is optimal

$$p(s_{1:T}, a_{1:T} \mid \mathcal{O}_{1:T}) \propto p(s_{1:T}, a_{1:T}) \exp\left(\sum_{t=1}^T r(s_t, a_t)\right)$$

A Derivation from Probabilistic Graphical Model

- Control via Variational Inference.
Recall the ELBO

$$\log p(O_{1:T}) \geq \mathbb{E}_{s_{1:T}, a_{1:T} \sim q} \left[\log \frac{p(O_{1:T}, s_{1:T}, a_{1:T})}{q(s_{1:T}, a_{1:T})} \right]$$

Because q is an arbitrary distribution, let

$$q(s_{1:T}, a_{1:T}) = p(s_1) \prod_{t=1}^T p(s_{t+1} | s_t, a_t) \pi(a_t | s_t)$$

then

$$\log p(O_{1:T}) \geq \sum_{t=1}^T \mathbb{E}_{s_t, a_t \sim q} [r(s_t, a_t) + \mathcal{H}(\pi(a_t | s_t))]$$

Preliminary 2: Soft Bellman Equation

- Soft Q function

$$Q_{\text{soft}}^{\pi}(s, a) = r(s, a) + \mathbb{E}_{s_t, a_t \sim \rho_{\pi}} \left[\sum_{t=1}^{\infty} \gamma^t (r(s_t, a_t) + \alpha H(\pi(\cdot | s_t))) \mid s_0 = s, a_0 = a \right]$$

Soft Bellman equation

$$Q_{\text{soft}}^{\pi}(s, a) = r(s, a) + \gamma \mathbb{E}_{\pi} [Q_{\text{soft}}^{\pi}(s_{t+1}, a_{t+1}) + \alpha \mathcal{H}(\pi(\cdot | s_{t+1})) \mid s_t = s, a_t = a]$$

- Energy Based Policy

Theorem

Given a policy π , define a new policy $\tilde{\pi}$ as

$$\tilde{\pi}(\cdot | s) \propto \exp(Q_{\text{soft}}^{\pi}(s, \cdot)), \quad \forall s$$

Assume that Q and $\int \exp(Q(s, a)) da$ are bounded. Then

$$Q_{\text{soft}}^{\tilde{\pi}(s, a)} \geq Q_{\text{soft}}^{\pi}, \quad \forall s, a$$

Soft Bellman Optimality Equation

The optimal soft Q function $Q_{\text{soft}}^*(s, a) = Q_{\text{soft}}^{\pi_{\text{MaxEnt}}^*}(s, a)$ satisfies

$$\pi_{\text{MaxEnt}}^* \propto \exp\left(\frac{1}{\alpha} Q_{\text{soft}}^*\right)$$

Then we get the soft Bellman optimality equation

$$Q_{\text{soft}}^*(s_t, a_t) = r_t + \gamma \mathbb{E}_{\pi} \left[\alpha \cdot \text{softmax} \frac{1}{\alpha} Q_{\text{soft}}^*(s_{t+1}, \cdot) \right]$$

where

$$\text{softmax} Q(s, \cdot) = \log \int_{\mathcal{A}} \exp(Q(s, a)) da$$

Soft Q-learning (SQL)²

- We aim to minimize the soft Bellman error

$$\min_Q |\mathcal{T}Q - Q|$$

where

$$\mathcal{T}Q(s, a) \triangleq r(s, a) + \gamma \mathbb{E}_{s' \sim p_s} \left[\log \int \exp Q(s', a') da' \right]$$

- We solve it via *soft policy iteration*

$$Q_{\text{soft}}^{(k+1)} \rightarrow \mathcal{T}Q_{\text{soft}}^{(k)}$$

- Theoretic guarantee: $Q_{\text{soft}}^{(k)}$ convergence to a local minima.

²Tuomas Haarnoja et al. "Reinforcement learning with deep energy-based policies". In: *International conference on machine learning*. PMLR, 2017, pp. 1352–1361.

Soft Policy Iteration in SQL

Similar to standard Q learning, we iteratively do the following three steps:

- Collect experience (s_t, a_t, s_{t+1}, r_t) from the environment with π
- Soft Policy Evaluation (update Q)

$$\hat{Q}(s_t, a_t) = r_t + \gamma \cdot \alpha \cdot \text{softmax} \frac{1}{\alpha} Q_{\text{soft}}(s_{t+1}, \cdot)$$

$$J_Q = \mathbb{E}_{(s_t, a_t, r_t, s_{t+1}) \sim \mathcal{D}} \frac{1}{2} (\hat{Q}(s_t, a_t) - Q(s_t, a_t))^2$$

update Q with the gradient of J_Q

- Soft Policy Improvement (update π)

$$\pi(\cdot | s) \propto \exp\left(\frac{1}{\alpha} Q(s, \cdot)\right)$$

Soft Bellman Equation

- Soft value function

$$V_{\text{soft}}^{\pi}(s) = \mathbb{E}_{s_t, a_t \sim \rho_{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \alpha H(\pi(\cdot | s_t))) \mid s_0 = s \right]$$

Soft Bellman equation

$$\begin{aligned} Q_{\text{soft}}^{\pi}(s, a) &= r(s, a) + \gamma \mathbb{E}_{\pi} [(Q_{\text{soft}}^{\pi}(s_{t+1}, a_{t+1}) + \alpha \mathcal{H}(\pi(\cdot | s_{t+1}))) \mid s_t = s, a_t = a] \\ &= r(s, a) + \gamma \mathbb{E}_{\pi} [V_{\text{soft}}^{\pi}(s_{t+1}) \mid s_t = s, a_t = a] \end{aligned}$$

$$\begin{aligned} V_{\text{soft}}^{\pi}(s) &= \mathbb{E}_{\pi} [Q_{\text{soft}}^{\pi}(s_t, a_t) \mid s_t = s] + \alpha \mathcal{H}(\pi(\cdot | s)) \\ &= \mathbb{E}_{\pi} [Q_{\text{soft}}^{\pi}(s_t, a_t) - \alpha \log \pi(a_t | s_t) \mid s_t = s] \end{aligned}$$

Soft Actor-Critic (SAC)³

- We aim to minimize the soft Bellman error

$$\min_Q |\mathcal{T}Q - Q|$$

where

$$\mathcal{T}Q(s_t, a_t) \triangleq r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim p} [V(s_{t+1})]$$

where

$$V(s_t) = \mathbb{E}_{a_t \sim \pi} [Q(s_t, a_t) - \log \pi(a_t | s_t)]$$

- We solve it via *soft policy iteration*

$$Q_{\text{soft}}^{(k+1)} \rightarrow \mathcal{T}Q_{\text{soft}}^{(k)}$$

- Theoretic guarantee: $Q_{\text{soft}}^{(k)}$ convergence to a local minima.

³Tomas Haarnoja et al. "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor". In: *International conference on machine learning*. PMLR, 2018, pp. 1861–1870.

Soft Policy Iteration in SAC

- Collect experience
- Soft Policy Evaluation

$$\hat{Q}(s_t, a_t) = r_t + \gamma V(s_{t+1})$$

$$\hat{V}(s_t, a_t) = Q(s_t, a_t) - \alpha \log \pi(a_t | s_t)$$

$$J_Q = \frac{1}{2} \mathbb{E}_{(s_t, a_t, r_t, s_{t+1}) \sim \mathcal{D}} (\hat{Q}(s_t, a_t) - Q(s_t, a_t))^2$$

$$J_V = \frac{1}{2} \mathbb{E}_{(s_t, a_t, r_t, s_{t+1}) \sim \mathcal{D}} (\hat{V}(s_t, a_t) - V(s_t, a_t))^2$$

- Soft Policy Improvement

$$\pi(\cdot | s) \propto \exp\left(\frac{1}{\alpha} Q(s, \cdot)\right)$$

Ablation: Stochastic vs. Deterministic Policy

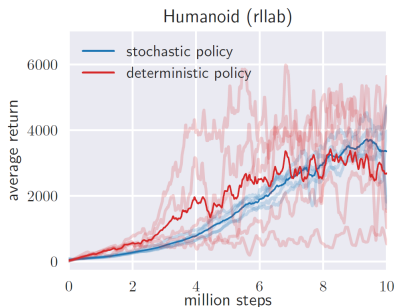
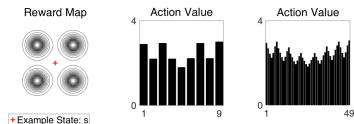
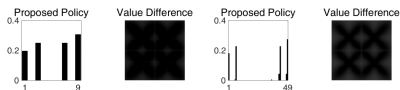
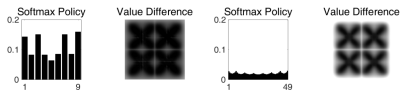


Figure 2. Comparison of SAC (blue) and a deterministic variant of SAC (red) in terms of the stability of individual random seeds on the Humanoid (rllab) benchmark. The comparison indicates that stochasticity can stabilize training as the variability between the seeds becomes much higher with a deterministic policy.

- 1 Introduction
- 2 Soft MDP
- 3 Learn a Sparse Policy**
- 4 Application: Mobile Health

A Toy Example⁴(a) Reward map and action values at state s .

(b) Proposed policy model and value differences (darker is better).



(c) Softmax policy model and value differences (darker is better).

Figure 4: A 2-dimensional multi-objective environment with point mass dynamics. The state is a location and the action is a velocity bounded with $[-3, 3] \times [-3, 3]$. The action space is discretized into two levels: 9 (low resolution) and 49 (high resolution).

From *Soft* to *Sparse*

- Soft: assigns a non-negligible probability mass to all actions
- Sparse: weeds out suboptimal actions and maintains near optimal actions

We want a sparse policy instead of a soft policy for two reasons:

- Safety: the policy gives bad actions with zero probability
- Performance

Regularized MDP⁵

- Regularized MDP

$$\pi_{\lambda}^* = \arg \max_{\pi} \mathbb{E}_{\pi} \sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \lambda \phi(\pi(a_t | s_t)))$$

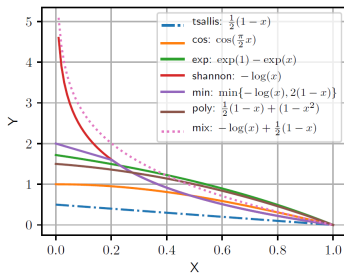
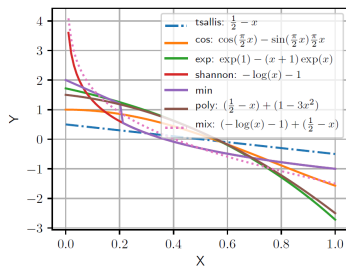
where $\phi(x) \in [0, 1]$ is a regularization term.

- Assumption for the regularization term

- ▶ Monotonicity: $\phi(x)$ is non-increasing
- ▶ Non-negativity: $\phi(1) = 0$
- ▶ Differentiability: $\phi(x)$ is differentiable;
- ▶ **Expected Concavity**: $f_{\phi}(x) = x\phi(x)$ is strictly concave

⁵Wenhao Yang, Xiang Li, and Zhihua Zhang. "A regularized approach to sparse optimal policy in reinforcement learning". In: *Advances in Neural Information Processing Systems* 32 (2019).

Examples of Regularization Terms

(a) Different ϕ 's(b) f'_ϕ for different ϕ 's

Bellman Optimality Equation

- Value function

$$V_{\lambda}^{\pi}(s) = \mathbb{E}_{\pi} \sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \lambda \phi(\pi(a_t | s_t)))$$

- Q function and Bellman equation

$$\begin{aligned} Q_{\lambda}^{\pi}(s, a) &= r(s, a) + \gamma \mathbb{E} V_{\lambda}^{\pi}(s') \\ &= r(s, a) + \gamma \mathbb{E}_{\pi} [Q_{\lambda}^{\pi}(s', a') + \lambda \phi(\pi(a' | s'))] \end{aligned}$$

- Bellman Optimality Equation

$$\begin{aligned} Q_{\lambda}^*(s, a) &= r(s, a) + \gamma \mathbb{E} V_{\lambda}^*(s') \\ \pi^*(a | s) &= \max \left\{ g_{\phi} \left(\frac{\mu_{\lambda}^*(s) - Q_{\lambda}^*(s, a)}{\lambda} \right), 0 \right\} \\ V_{\lambda}^*(s) &= \mu_{\lambda}^* - \lambda \sum_a \pi_{\lambda}^*(a | s)^2 \phi'(\pi_{\lambda}^*(a | s)) \end{aligned}$$

where $g_{\phi} = (f'_{\phi})^{-1}$

Definition and Property of Sparsity

Definition

A given policy π is called δ -sparse if it satisfies:

$$\frac{|\{(s, a) \in \mathcal{S} \times \mathcal{A} \mid \pi(a \mid s) \neq 0\}|}{|\mathcal{S}||\mathcal{A}|} \leq \delta$$

If $\pi(a \mid s) > 0$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, we call it has no sparsity.

Theorem

If $\lim_{x \rightarrow 0^+} f'_\phi(x) = \infty$, the optimal policy has no sparsity.

Remark

- Proof is directly from the bellman equation
- Entropy regularization has no sparsity
- If $f'_\phi(0) < \infty$, we can control the sparsity with λ

Performance Error between V_λ^* and V^*

Under some assumptions, we have

$$\|V_\lambda^* - V^*\|_\infty \leq \frac{\lambda}{1-\gamma} \phi\left(\frac{1}{|\mathcal{A}|}\right)$$

Regularized Actor-Critic (RAC)

- Collect experience
- Policy evaluation

$$\hat{Q}(s_t, a_t) = r_t + \gamma(Q(s_{t+1}, a_{t+1}) + \lambda\phi(\pi(a_{t+1} | s_{t+1})))$$

$$J_Q = \frac{1}{2} \mathbb{E}_{\mathcal{D}}(\hat{Q} - Q)$$

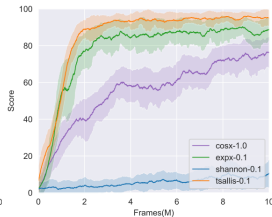
- Policy improvement

$$J_{\pi} = -\mathbb{E}_{\mathcal{D}}[Q + \lambda\phi(\pi)]$$

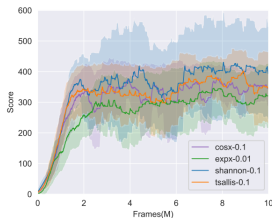
Experience on Atari games



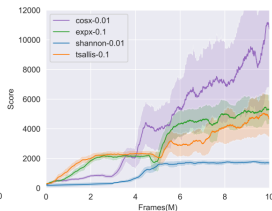
(a) Alien



(b) Boxing



(c) Breakout



(d) Seaquest

- 1 Introduction
- 2 Soft MDP
- 3 Learn a Sparse Policy
- 4 Application: Mobile Health**

Background: mobile health (mHealth) technology

mHealth (mobile health) is a general term for the use of mobile phones and other wireless technology in medical care.

- Usages
 - ▶ educations about preventive healthcare services
 - ▶ disease surveillance
 - ▶ **treatment support**
 - ▶ epidemic outbreak tracking
 - ▶ **chronic disease management**
- Advantages: convenient
- Disadvantages
 - ▶ data breach
 - ▶ information may not be accurate
 - ▶ **safety**

Treatment Support⁶

We aim to monitor individuals' health statuses and deliver just-in-time personalized interventions. For example, monitor a diabetic's blood sugar, heart rate and other physiological indicators and decide the dose of insulin injection.

- large numbers of intervention options
- unclear what alternatives can be used (when meet temporary medication shortages)

⁶Wenzhuo Zhou, Ruoqing Zhu, and Annie Qu. "Estimating optimal infinite horizon dynamic treatment regimes via pt-learning". In: *Journal of the American Statistical Association* just-accepted (2022), pp. 1–40.

Methodology and Theory

- choose $\phi(x) = -\frac{1}{2}(x - 1)$
- on-line \rightarrow off-line, intuitively

$$\min_{\pi_\lambda} \|\mathcal{T}Q - Q\| \rightarrow \min_{\pi_\lambda} \sum_{i=1}^n \|\mathcal{T}Q(s^{(i)}, a^{(i)}) - Q(s^{(i)}, a^{(i)})\|$$

so we must consider the sample complexity

$$\|V_{\lambda}^{\hat{\pi}_\lambda} - V^*\|_{L^2}^2 = \mathcal{O}(n^{-1/2}) + \mathcal{O}(\lambda^2 |\mathcal{A}|)$$

where $\hat{\pi}_\lambda$ is the empirical risk minimizer and n is the number of samples.

Comments of Statistics Professionals: This theoretic guarantee is very important in statistics, since we can not evaluate our policy in the real world.

Simulation

- State: blood glucose, Adiponectin, blood pressure
- Action: insulin injections (Yes/No), physical activity (No/Moderate/Strong) and dietary intake (Yes/No)
- Reward: linear combination of the three states
- Transition: a given model for generating data (unknown to the policy)

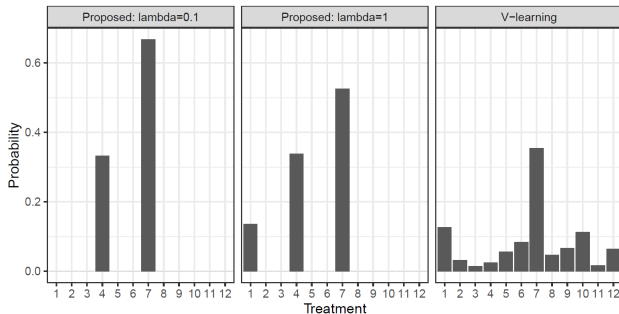


Figure 5: The estimated optimal policy distribution for one patient at a specific time point. The 7th treatment is the optimal treatment, and the 4th treatment is a near-optimal treatment; other treatments are sub-optimal or non-optimal.