$$TD(0): \quad V(s_t) \leftarrow V(s_t) + \alpha(r + \gamma V(s_{t+1}) - V(s_t))$$

- 1-step bootstrap vs. multi-step ✓ ←
- on-policy data vs. off-policy data.
- tabular representation vs. function approximation.
- value prediction. vs. control

# Value Prediction
# with Function Approximation
Reading: Algs for RL (Szepesvári), Sec 3.2

$$s_1, a_1, r_1, s_2. \cdots$$

$$V^\pi.$$

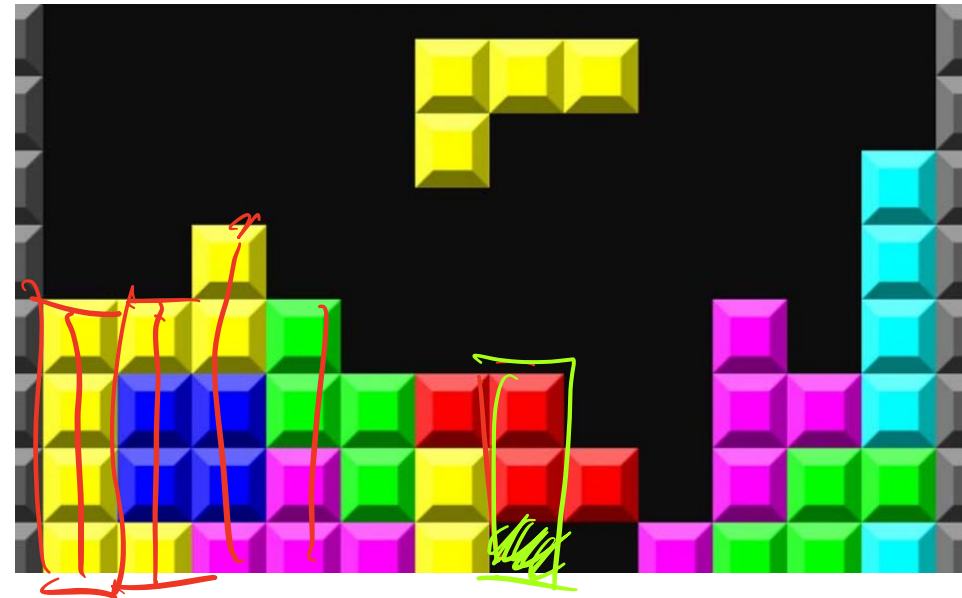# Generalization for value prediction

- Major limitation of tabular RL: does not scale to large state space
  - most methods require that we run into the same state multiple times
  - when the state space is large, you might not see the same state even twice!
  - In other words: sample complexity scales with |S|
  - need generalization
- For value prediction problem, generalization requires that we have some prior knowledge about the form the value function
  - linear function approximation: design features $\phi(s) \in \mathbb{R}^d$ ("featurizing states"), and approximate $V^\pi(s) \approx \theta^\top \phi(s)$
  - only unknown: $\theta$. $d$ unknowns vs |S| unknowns!

# Example: linear function approximation for tetris

- An example featurization:
  - let the height of the pile in i-th column be the i-th feature
  - dimensionality of feature = #columns
  - (probably doesn't work; just an example)
- Feature engineering requires a lot of prior knowledge, domain insights, and trial and error, just as in supervised learning!

$$V^\pi \approx \phi(s)^\top \theta + b$$

$$= [\phi(s) \ 1] \begin{bmatrix} \theta \\ b \end{bmatrix}$$

# Monte-Carlo Value Prediction

Handwritten top: $s_1^{(i)}, a_1^{(i)}, r_1^{(i)}, \dots, s_H^{(i)}, a_H^{(i)}, r_H^{(i)}$

Handwritten: $p(s) > 0 \ \forall s.$

- Draw a starting state $s_i$ from the exploratory initial distribution, roll out a trajectory using $\pi$ from $s_i$, and let $G_i$ be the (random) discounted return

Handwritten: $= \sum_{t=1}^{H} \gamma^{t-1} r_t^{(i)}$

Handwritten: $V^\pi(s) = \mathbb{E}[G \mid s]$

- Collect $\{(s_i, G_i)\}$ pairs

- Least square regression: $\min_\theta \frac{1}{n} \sum_{i=1}^{n} (\theta^\top \phi(s_i) - G_i)^2$

- Why this works?

Handwritten: $V^\pi(s) = \mathbb{E}[G \mid s] = \underset{f: S \to \mathbb{R}}{\arg\min} \ \mathbb{E}[(f(s) - G)^2]$
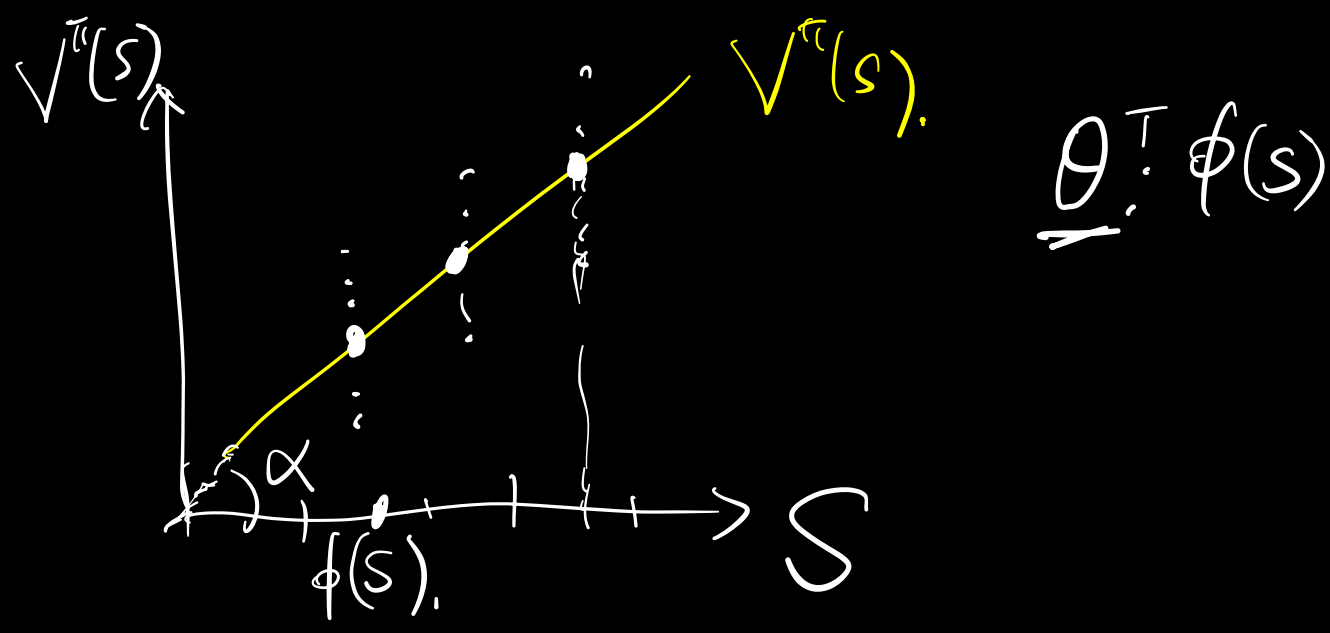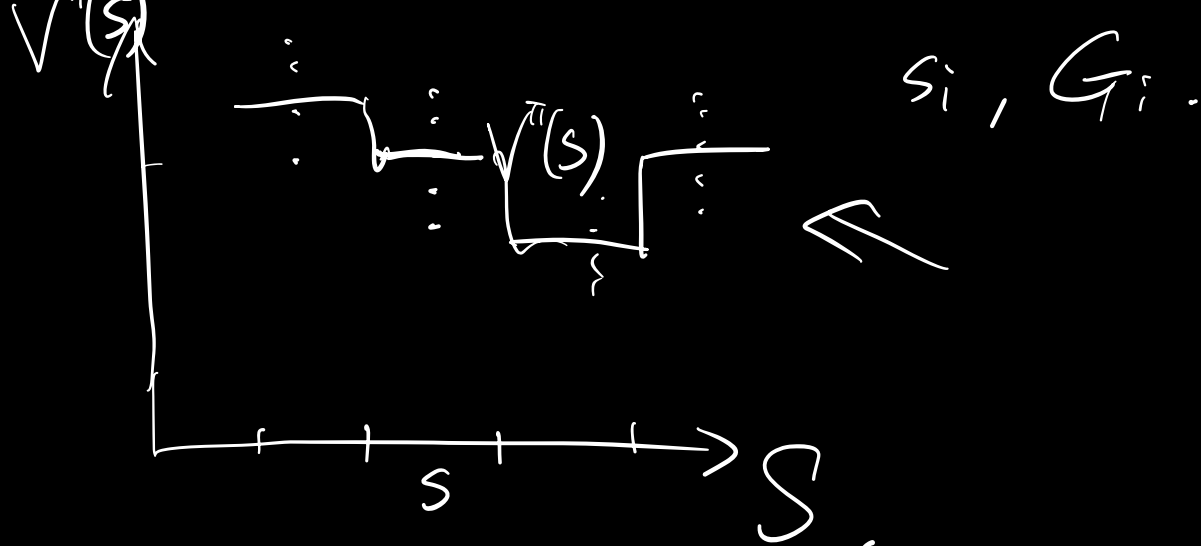
  - Assume $\{(s_i, G_i)\}$ are i.i.d., let $(s, G)$ be variables equal in distribution

  - The expected version of the objective: $\min_\theta \mathbb{E}_{s, G}[(\theta^\top \phi(s) - G)^2]$

  - If we do not restrict ourselves to linear functions, the function that minimizes this objective is $s \mapsto \mathbb{E}[G \mid s] \ \ (= V^\pi(s))$!

  - If true $V^\pi(s)$ happens to take linear form, the regression will find it in the limit (of infinite data)

  - Finite sample regime: bias & variance trade-off

4

$$\mathbb{E}\left[\left(f(x)-Y\right)^2\right] = \mathbb{E}\left[\left(f(x)-\mathbb{E}[Y|x]\right)^2\right] \leftarrow$$

$$+ \mathbb{E}\left[\left(\mathbb{E}[Y|x]-Y\right)^2\right].$$

$$S \quad \frac{1}{n}\sum_{i=1}^{n}\left(G_i\right).$$

$$V(S_t) \leftarrow V(S_t) + \alpha\left(G_t - V(S_t)\right).$$

$$\mathbb{R}^S \ni V \leftarrow V + \alpha\left(G_t - V(S_t)\right)\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad t\text{-th}$$

$$\theta.$$

$$\begin{bmatrix} 0, 0 \dots 1 0 0 \end{bmatrix} \begin{bmatrix} V(S') \\ V(S^2) \\ \vdots \\ V(S) \end{bmatrix} \quad \phi(S_t).$$

$$\theta \leftarrow \theta + \alpha\left(G_t - \phi(S_t)^\top\theta\right)\phi(S_t).$$

$V^\pi(s)$

$V^\pi(s)$

$s$

$S$

$s_i, G_i.$

$V^\pi(s)$

$\theta^T \phi(s)$

$V^\pi(s)$

$\alpha$

$\phi(s).$

$S$

# Monte-Carlo Value Prediction

$$\underset{V_\theta \in \mathcal{F}}{\arg\min} \; \mathbb{E}[(V_\theta(s) - G)^2]$$

- The same idea applies to non-linear value function approximation
- More generally & abstractly, think of function approximation as searching over a restricted **function space**, which is a set whose members are functions that map states to real values.
- Function space of linear value function approximation: $\mathcal{F} = \{V_\theta : \theta \in \mathbb{R}^d\}$, where $V_\theta(s) = \theta^\top \phi(s)$
  - typically only a small subset of all possible functions
  - Using "all possible functions" = tabular!
  - Equivalently, tabular MC value prediction can be recovered by choosing $\phi$ as the identity features $\phi(s) = \{\mathbb{I}[s = s']\}_{s' \in S}$
- $\min_{V_\theta \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} (V_\theta(s_i) - G_i)^2$
- Plug in any function approximator of your choice
- SGD: uniformly sample i and $\theta \leftarrow \theta - \alpha \cdot (V_\theta(s_i) - G_i) \cdot \nabla V_\theta(s_i)$

$$\nabla_\theta (\phi(s)^\top \cdot \theta) = 0$$

$$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ S \rightarrow 1 \\ \vdots \\ 0 \end{bmatrix}$$

$$V(s_t) \leftarrow V(s_t) + \alpha \left( V_t + \gamma V(s_{t+1}) - V(s_t) \right)$$

$$\boxed{V_{k+1} \leftarrow T^\pi V_k.}$$

$$s \quad \{(r_i, s_i')\}_{i=i}^n$$

$$V_{k+1}(s) = \mathbb{E}_\pi \left[ r + \gamma V_k(s') \quad s \right].$$

$$= \underset{f: S \to \mathbb{R}}{\text{argmin}} \; \mathbb{E}_\pi \left[ \left( \underbrace{f(s)}_{\text{input}} - \underbrace{(r + \gamma V_f(s'))}_{\text{label}} \right)^2 \right].$$

$$\approx \underset{V_\theta \in \mathcal{F}}{\text{argmin}} \quad \frac{1}{n} \sum_{i=1}^n \left( V_\theta(s) - r_i - \gamma V_k(s_i') \right)^2.$$

$$\text{SGD:} \quad \theta \leftarrow \theta + \alpha \left( V_\theta(s) - r_i - \gamma V_k(s_i') \right) \nabla V_\theta(s)$$

$$s_1, a_1, r_1, \quad \cdots \cdots \cdots, \quad s_t, a_t, V_\theta, \quad \cdots$$

$$\theta \leftarrow \theta + \alpha \left( V_\theta(s_t) - r_t - \gamma V_\theta(s_{t+1}) \right) \nabla V_\theta(s_t).$$

# TD(0) with function approximation

- tabular: $V(s_t) \leftarrow V(s_t) + \alpha(r_t + \gamma V(s_{t+1}) - V(s_t))$

- When we update $V(s_t)$, the target is $r_t + \gamma V(s_{t+1})$

- Batch version of the algorithm: one Bellman update can be approximated (using all data) as

$$V_{k+1} \leftarrow \arg\min_{V_\theta \in \mathcal{F}} \frac{1}{T} \sum_{t=1}^{T} (V_\theta(s_t) - r - \gamma V_k(s_{t+1}))^2$$

- SGD + "no-wait": $\theta \leftarrow \theta - \alpha \cdot (V_\theta(s_t) - r - \gamma V_\theta(s_{t+1})) \cdot \nabla V_\theta(s_t)$

- When using linear function approximation $V_\theta(s) = \phi(s)^\top \theta$, we have
  $\theta \leftarrow \theta - \alpha \cdot (\phi(s_t)^\top \theta - r - \gamma \phi(s_{t+1})^\top \theta) \cdot \phi(s_t)$

- When using chain rule, we only take gradient on $V_\theta(s_t)$ and ignore $V_\theta(s_{t+1})$; the latter is treated as a constant (it plays the role of $V_k$)

$$\underline{\theta^T \phi(s)} \approx \boxed{R(s,\pi) + r \, \underset{s'|s,\pi}{\mathbb{E}} \left[ \theta^T \phi(s') \right].}$$

$$\theta \text{ s.t. } \phi(s)^T \theta \approx V^\pi(s).$$

$$s, \quad a \sim \pi, \quad r, \quad s'.$$

$$\mathbb{E} \left[ \left( \theta^T \phi(s) - r - \gamma \; \theta^T \phi(s') \right)^2 \right]$$

$$\theta \leftarrow \theta + \underset{=}{\alpha} \cdot \left( \theta^T \phi(s) - r - \gamma \, \theta^T \phi(s') \right)$$
$$\left( \phi(s) - \gamma \phi(s') \right).$$

# What if…? (not required)

- What happens if we also differentiate $V_\theta(s_{t+1})$?
- This corresponds to $\arg\min_{V_\theta \in \mathcal{F}} \sum_{(s,r,s')} (V_\theta(s) - r - \gamma V_\theta(s'))^2$
  - no iteration anymore; a clean optimization objective
  - (most RL algorithms with bootstrapped target do not have a fixed optimization objective; objective changes over time)
- Assume for simplicity that, each data point is generated by (1) sampling $s$ i.i.d. from some exploratory distribution, and (2) generating $r$ and $s'$ conditioned on $(s, \pi(s))$
- Replacing empirical objective with the population version, the objective becomes $\mathbb{E}_{s,r,s'}[(V_\theta(s) - r - \gamma V_\theta(s'))^2]$
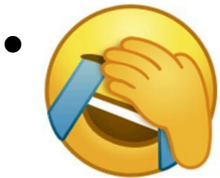
# What if…? (not required)

- $\mathbb{E}_{s,r,s'}[(V_\theta(s) - r - \gamma V_\theta(s'))^2]$ can be decomposed into two terms
  - First term: $\mathbb{E}_s[(V_\theta(s) - (\mathscr{T}^\pi V_\theta)(s))^2]$
    - This is good! measures how much $V_\theta$ violates Bellman eq
    - A version of Bellman error $\|V - \mathscr{T}^\pi V\|$
  - Second term: $\gamma^2 \mathbb{E}_s[\mathrm{Var}_{s'|s,\pi(s)}[V_\theta(s')]]$
    - (assumes deterministic rewards)
    - This is bad! An additional term that penalizes functions that has large variance w.r.t. random state transitions
    - Special case: 0 when environment is deterministic
- So it's actually a sensible algorithm for deterministic environments, but doesn't work when stochasticity is significant

# Resolutions (not required)

- If we have a simulator…
  - For each s in data, draw another independent state transition
  - Minimize objective $\mathbb{E}[(V_\theta(s) - r - \gamma V_\theta(s'_A))(V_\theta(s) - r - \gamma V_\theta(s'_B))]$
  - "Double sampling" and Baird's residual algorithm (Bellman residual minimization)
  - Exercise: do you need to double sample the reward if reward is stochastic?
  - The conditional variance term is eliminated by double sampling
- If we can only draw 1 next-state (as with any natural data generation process)…
  - Estimate the conditional variance term and subtract from the objective
  - A minimax formulation (not covered in this course)
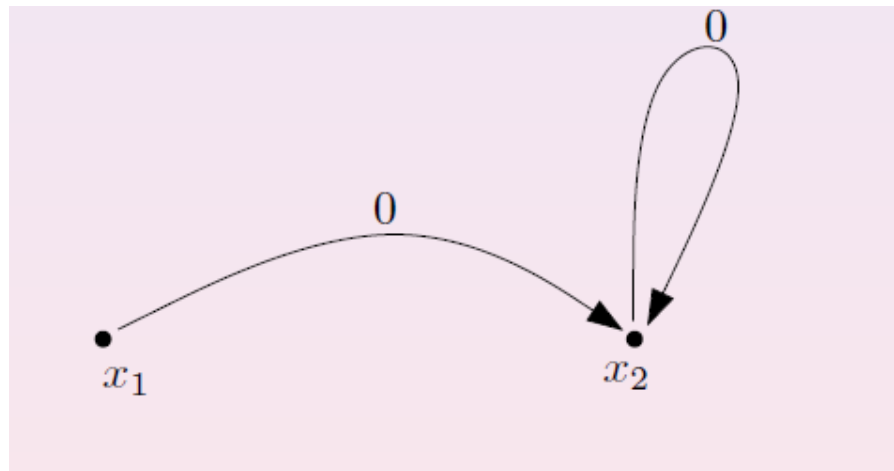  - For further readings, see 542 slides on FQI.

# Convergence?

- TD with function approximation can diverge in general
- Is it because of…
  - Randomness in SGD?
    - Nope. Even the batch version doesn't converge.
  - Sophisticated, non-linear func approx?
    - Nope. Even linear doesn't converge.
  - That our function class does not capture $V^\pi$?

    - Nope. Even if $V^\pi$ can be exactly represented in the function class ("realizable"), it still does not converge.

- 🤦

## 2.1 Counter-example for least-square regression [Tsitsiklis and van Roy, 1996]

An MDP with two states $x_1, x_2$, 1-d features for the two states: $f_{x_1} = 1, f_{x_2} = 2$. Linear Function approximation with $\tilde{V}_\theta(x) = \theta f_x$.



$$
\begin{aligned}
\theta_k \quad &:= \quad \arg\min_\theta \frac{1}{2}(\theta - \text{target}_1)^2 + (2\theta - \text{target}_2)^2 \\
&= \quad \arg\min_\theta \frac{1}{2}(\theta - \gamma\theta^{k-1} f_{x_2})^2 + (2\theta - \gamma\theta^{k-1} f_{x_2})^2 \\
&= \quad \arg\min_\theta \frac{1}{2}(\theta - \gamma 2\theta^{k-1})^2 + (2\theta - \gamma 2\theta^{k-1})^2
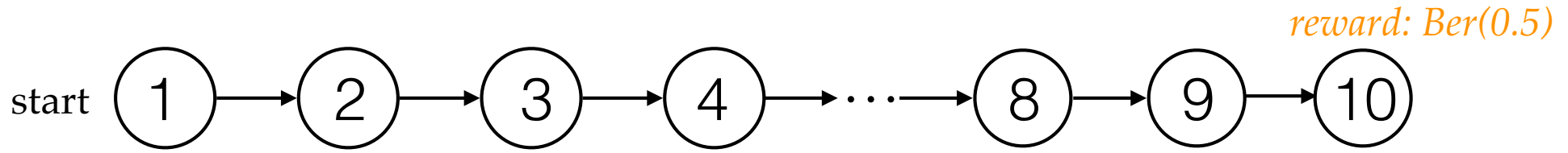\end{aligned}
$$

$$
(\theta - \gamma 2\theta^{k-1}) + 2(2\theta - \gamma 2\theta^{k-1}) = 0 \Rightarrow 5\theta = 6\gamma\theta^{k-1}
$$

$$
\theta_k = \frac{6}{5}\gamma\theta_{k-1}
$$

This diverges if $\gamma \geq 5/6$.

credit: course notes
from Shipra Agrawal

# A simple example (finite horizon, $\gamma=1$)

*reward: Ber(0.5)*



start $\rightarrow$ (1) $\rightarrow$ (2) $\rightarrow$ (3) $\rightarrow$ (4) $\rightarrow$ ... $\rightarrow$ (8) $\rightarrow$ (9) $\rightarrow$ (10)

Iter #1:   Data: (⑩, 1, *end*), ..., (⑩, 0, *end*) $\Rightarrow$   0.501

Iter #2:   Data: (⑨, 0, ⑩) $\Rightarrow$ (⑨, 0+0.501) $\Rightarrow$   0.501   0.501

...

Iter #10:  0.501   0.501   0.501   0.501   ...   0.501   0.501   0.501

- Dataset $D = \{(s, r, s')\}$ looks like:
  $\{(①, 0, ②), (②, 0, ③), ..., (⑩, 1, end), ..., (⑩, 0, end)\}$

12

# How things go wrong (w/ restricted class)

*reward: Ber(0.5)*

start  (1) → (2) → (3) → (4) → ... → (8) → (9) → (10)

Realizable

Function class

| 0.5 | × | 0.5 | × | 0.5 | × | ... | × | 0.5 | × | 0.5 | × | 0.5 |

1.012  0.756  0.628  ...  0.504  0.502  0.501

Iter #1:  Data: (⑩, 1, *end*), ..., (⑩, 0, *end*)  ⇨  0.501

Iter #2:  Data: (⑨, 0, ⑩)  ⇨  (⑨, 0+0.501)  ⇨  0.502  0.501

...

**!!!**

Iter #10: 1.012  0.756  0.628  ...  0.502  0.501

Example given in Dann et al'18

13

# Non-convergence

- Why things go wrong?
- Bellman update is a contraction, but here we have an additional projection step: $V_{k+1} \leftarrow \Pi_{\mathcal{F}} \mathcal{T}^{\pi} V_k$, projected Bellman update may NOT be a contraction (even with linear function approximation)
  - it is still a contraction in some special cases; will see
- In other words: in each iter, we solve a regression problem where the target function is $T^{\pi} V$, where $V$ can be arbitrary function in $\mathcal{F}$
- The fact that $V^{\pi} \in \mathcal{F}$ does not imply that $T^{\pi} V$ is in $\mathcal{F}$! We may do quite poorly in the regression problem, and the iteration does not mimic a Bellman update
- Why tabular is fine? $\mathcal{F}$ is fully expressive so $T^{\pi} V$ is always in $\mathcal{F}$.
  - Similarly for func approx, if we assume that $\mathcal{F}$ is closed under $T^{\pi}$, can prove some good properties of TD
- All alg based on bootstrapped targets suffer from this issue
  - Compare to the behavior of Monte-Carlo