

Fitted Q-Iteration

(most references can be found
on paper list for project topics)

Generalization for value-based batch RL

- We studied using abstractions to generalize in large state spaces
- Abstractions correspond to “histogram regression” in supervised learning—the most trivial form of generalization
- Can I use XXX for value-based RL?
 - Linear predictors?
 - Kernel machines?
 - Random forests?
 - Neural nets???
 - ...
- **What you really want:** *Reduction* of RL to supervised learning.

always min f $\mathbb{E}[(Y-f(x))^2]$.

(s, a, r, s')

Revisiting value iteration

- Recall the value iteration algorithm: $f_k \leftarrow \mathcal{T} f_{k-1}$

$(s, a), Y$

- where $(\mathcal{T} f)(s, a) = \mathbb{E}_{r \sim R(s, a), s' \sim P(\cdot | s, a)} [r + \gamma \max_{a'} f(s', a')]$

- i.e., $\mathcal{T} f_{k-1} = \mathbb{E} [r + \gamma \max_{a'} f_{k-1}(s', a') \mid s, a]$

- What we want: a function in the form of $\mathbb{E}[Y | X]$

$= \mathcal{T} f_{k-1}$

- $Y = r + \gamma \max_{a'} f_{k-1}(s', a')$, $X = (s, a)$

- How to obtain $\mathbb{E}[Y | X]$? Squared-loss regression!!!

- Fitted-Q Iteration [Ernst et al'05]

$f^* = \mathcal{T} f_{k-1}$

$$f_k = \arg \min_{f \in \mathcal{F}} \sum_{(s, a, r, s') \in D} \left(f(s, a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f_{k-1}(s', a') \right) \right)^2$$

- $F =$ all functions: FQI = VI in the estimated tabular model
- $F =$ all piece-wise const functions under abstraction ϕ : FQI = VI in the estimated abstract model

$(X_1, Y_1), (X_2, Y_2) \dots (X_n, Y_n)$

$$\boxed{(X, Y) \sim P.}$$

$$f^* = \operatorname{argmin}_{f: X \rightarrow \mathbb{R}} \mathbb{E}_P [(f(X) - Y)^2].$$

$$\forall x \in X \quad f(x) = \underline{f_x}.$$

$$\forall x. \operatorname{argmin}_{\underline{f_x}} \mathbb{E}_{Y \sim P(\cdot | x)} [(f_x - Y)^2].$$

$$\mathbb{E}_{Y \sim P(\cdot | X=x)} \left[(f_x - \mathbb{E}[Y | X=x])^2 + (\mathbb{E}[Y | X=x] - Y)^2 \right].$$

$$f_x = \mathbb{E}[Y | X=x].$$

$$\boxed{f^*(x) = \mathbb{E}[Y | X].}$$

Special case: MBRL (CE) with ϕ

$$D_{x,a} = \bigcup_{s \in \phi^{-1}(x)} D_{s,a}$$

- Algorithm: estimate \widehat{M}_ϕ , and do planning

$$\widehat{R}_\phi(x, a) = \frac{1}{|D_{x,a}|} \sum_{(r,s') \in D_{x,a}} r, \quad \widehat{P}_\phi(x, a) = \frac{1}{|D_{x,a}|} \sum_{(r,s') \in D_{x,a}} \mathbf{e}_{\phi(s')}$$

- Use Value Iteration as the planning algorithm:

- Initialize g_0 as any function in $\mathbb{R}^{S_\phi \times \mathcal{A}}$

- $g_t \leftarrow \mathcal{T}_{\widehat{M}_\phi} g_{t-1}$. That is, for each $x \in S_\phi$, $a \in A$:

$$g_t(x, a) = \widehat{R}_\phi(x, a) + \gamma \langle \widehat{P}_\phi(x, a), V_{g_{t-1}} \rangle$$

$$= \frac{1}{|D_{x,a}|} \sum_{(r,s') \in D_{x,a}} (r + \gamma \langle \mathbf{e}_{\phi(s')}, V_{g_{t-1}} \rangle)$$

$$= \frac{1}{|D_{x,a}|} \sum_{(r,s') \in D_{x,a}} (r + \gamma V_{g_{t-1}}(\phi(s')))$$

$\rightarrow \max_{a'} g_{t-1}(x', a')$

Rewrite in the original S

- Rewrite the algorithm so that $f_t = [g_t]_M$
- Define $\mathcal{F}^\phi \subset \mathbb{R}^{|S \times A|}$ as the space of all functions over $S \times A$ that are piece-wise constant under ϕ with value in $[0, V_{\max}]$
- Initialize f_0 as any function in F^ϕ
- For each $s \in S, a \in A$: essentially $f_t \leftarrow \mathcal{T}_{\widehat{M}'_\phi} f_{t-1}$

$$g_0, g_1, g_2, g_3, \dots \in \mathbb{R}^{|S \times A|}$$

$$\begin{aligned} f_t(s, a) &= \widehat{R}_\phi(\phi(s), a) + \gamma \langle \widehat{P}_\phi(\phi(s), a), [V_{f_{t-1}}]_\phi \rangle \\ &= \frac{1}{|D_{\phi(s), a}|} \sum_{(r, s') \in D_{\phi(s), a}} (r + \gamma \langle \mathbf{e}_{\phi(s')}, [V_{f_{t-1}}]_\phi \rangle) \\ &= \frac{1}{|D_{\phi(s), a}|} \sum_{(r, s') \in D_{\phi(s), a}} (r + \underbrace{\gamma V_{f_{t-1}}(s')}) \end{aligned}$$

“Empirical Bellman update”
(based on 1 data point)

$$\begin{aligned} g_t(x, a) &= \widehat{R}_\phi(x, a) + \gamma \langle \widehat{P}_\phi(x, a), V_{g_{t-1}} \rangle \\ &= \frac{1}{|D_{x, a}|} \sum_{(r, s') \in D_{x, a}} (r + \gamma \langle \mathbf{e}_{\phi(s')}, V_{g_{t-1}} \rangle) \\ &= \frac{1}{|D_{x, a}|} \sum_{(r, s') \in D_{x, a}} (r + \gamma V_{g_{t-1}}(\phi(s'))) \end{aligned}$$

$$f_t(s, a) = \frac{1}{|D_{\phi(s), a}|} \sum_{(r, s') \in D_{\phi(s), a}} \left(r + \gamma \max_{a' \in \mathcal{A}} f_{t-1}(s', a') \right)$$

Alternative interpretation of the above step

- Dataset $D = \{(s, a, r, s')\}$
- Apply emp. Bellman up. to f_{t-1} based on each data point:

$$\left\{ \left((s, a), \left(r + \gamma \max_{a' \in \mathcal{A}} f_{t-1}(s', a') \right) \right) \right\}$$

- What does it mean to take average over $D_{\phi(s), a}$?
 - Recall: average minimizes mean squared error (MSE)
 - *Projection* onto F^ϕ ! (think of functions over D)

$$f_t = \arg \min_{f \in \mathcal{F}^\phi} \sum_{(s, a, r, s') \in D} \left(f(s, a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f_{t-1}(s', a') \right) \right)^2$$

- ... which is, solving a SL regression problem with histogram regression F^ϕ

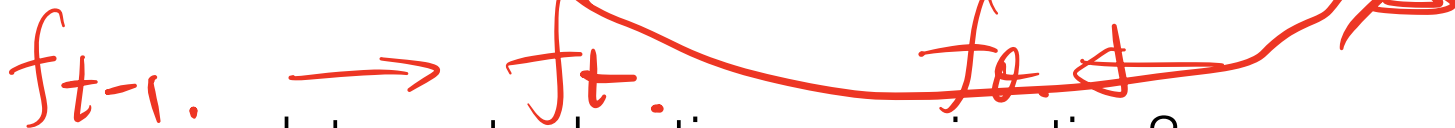
Fitted Q-Iteration (FQI): $f_t = \arg \min_{f \in \mathcal{F}} \sum_{(s,a,r,s') \in D} \left(f(s,a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f_{t-1}(s',a') \right) \right)^2$
 [Ernst et al'05]; see also [Gordon'95]



We simplified a “regression algorithm” to its corresponding function space F

- Empirical Risk Minimization (ERM); assume optimization is exact; does not consider regularization, etc.
- Will also assume finite (but exponentially large) F
 - continuous spaces are often handled by discretization in SLT (e.g., growth function, covering number)
 - methods like regression trees have dynamic function spaces (and often need SRM); not accommodated
- A minimal but (hopefully) insightful simplification of supervised learning

Fitted Q-Iteration (FQI): $f_t = \arg \min_{f \in \mathcal{F}} \sum_{(s,a,r,s') \in D} \left(f(s,a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f_{t-1}(s',a') \right) \right)^2$
 [Ernst et al'05]; see also [Gordon'95]



Asynchronous update + stochastic approximation?

- Assume parameterized & differentiable function: $\mathcal{F} = \{f_\theta : \theta \in \Theta\}$
- Online regression: randomly pick a data point and do a stochastic gradient update:

Treat as constant; don't pass gradient

$$\theta \leftarrow \theta - \frac{\alpha}{2} \cdot \nabla_\theta \left(f_\theta(s,a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f_\theta(s',a') \right) \right)^2$$

$$= \theta - \alpha \left(f_\theta(s,a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f_\theta(s',a') \right) \right) \nabla_\theta f_\theta(s,a)$$

Handwritten red annotations: Underlines under $\frac{\alpha}{2}$, $f_\theta(s,a)$, and $\nabla_\theta f_\theta(s,a)$. A red arrow points to the ∇_θ operator. A red scribble is next to the squared term.

- If f_θ is the tabular function, it's (tabular) Q-learning
- If f_θ is a neural net, it's (almost) DQN (Mnih et al.'15)
 - Using a target network is even more similar to FQI

Fitted Q-Iteration (FQI): $f_t = \arg \min_{f \in \mathcal{F}} \sum_{(s,a,r,s') \in D} \left(f(s,a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f_{t-1}(s',a') \right) \right)^2$
 [Ernst et al'05]; see also [Gordon'95]

The argmin step plays two roles:

1. Denoise the emp update $r + \gamma V_f(s')$ to $(\mathcal{T}f)(s, a)$ (w/ inf data)
 - This happens even in tabular setting
2. $\mathcal{T}f$ may not have a succinct representation => find the closest approximation in F (i.e., projection)
 - Denote Π_F as the projection. Dependence on weights over state-action pairs omitted—determined by data distribution
 - With infinite data, FQI becomes: $f_t \leftarrow \Pi_F \mathcal{T}f_{t-1}$

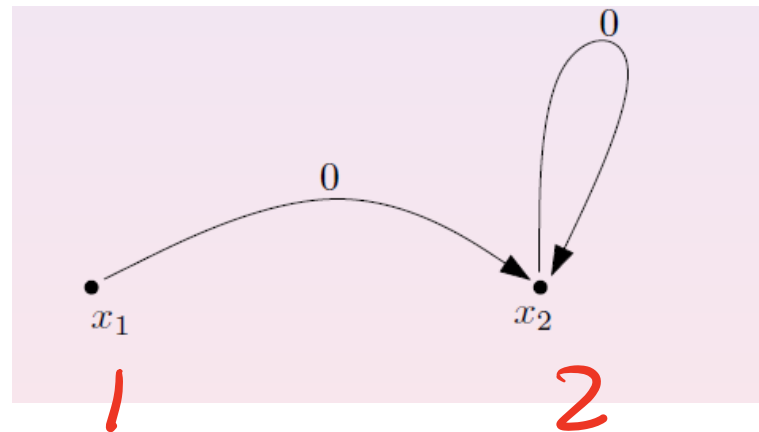
Convergence and Stability

- With infinite data, Q^* is a fixed point (as long as $Q^* \in F$)
 - $Q^* \in F$ is called (Q^* -) “*realizability*”
- CE w/ Q^* -irrelevant ϕ is a special case of FQI—convergence guaranteed
- Doesn't hold in general: FQI may diverge under $Q^* \in F$, even with
 - Infinite data
 - Fully exploratory data
 - Linear function class F
 - MDP has no actions (just policy evaluation)

$$(s, a) \sim \mu, r, s$$
$$\mu(s, a) > 0.$$

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 $\tilde{V}_\theta(x) = \theta f_x$.



$$V(x_2) = \phi(x_2) \cdot \theta$$

$$\downarrow$$

$$\frac{1}{2}$$

$$V(x_1)$$

$$= \phi(x_1) \cdot \theta.$$

credit: course notes from Shipra Agrawal

$$\theta_k := \arg \min_{\theta} \frac{1}{2} (\theta - \text{target}_1)^2 + \frac{1}{2} (2\theta - \text{target}_2)^2$$

$$= \arg \min_{\theta} \frac{1}{2} (\theta - \gamma \theta^{k-1} f_{x_2})^2 + \frac{1}{2} (2\theta - \gamma \theta^{k-1} f_{x_2})^2$$

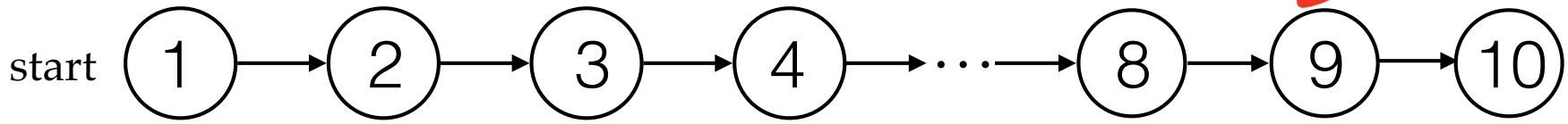
$$= \arg \min_{\theta} \frac{1}{2} (\theta - \gamma 2\theta^{k-1})^2 + \frac{1}{2} (2\theta - \gamma 2\theta^{k-1})^2$$

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

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

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

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



FQI
Iter #1:

Data: (⑩, 1, end), ..., (⑩, 0, end) ⇒

$$f_1(10) = 0.501$$

Iter #2:

Data: (⑨, 0, ⑩) ⇒ (⑨, 0+0.501) ⇒

0.501 0.501

$$r + \underbrace{f_{k-1}(s')}$$

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

- Dataset $D = \{(s, r, s')\}$ looks like (action omitted):
 $\{(\textcircled{1}, 0, \textcircled{2}), (\textcircled{2}, 0, \textcircled{3}), \dots, (\textcircled{10}, 1, \text{end}), \dots, (\textcircled{10}, 0, \text{end})\}$

Intuition for the instability

- Standard VI: $f_t \leftarrow \mathcal{T} f_{k-1}$
- FQI keeps things tractable by: $f_t \leftarrow \Pi_{\mathcal{F}}(\mathcal{T} f_{k-1})$
 - Π_F can destroy contraction of \mathcal{T} !
 - Preserved only in special cases (e.g., Q^* -irrelevant ϕ)
- A sufficient condition that fixes the issue

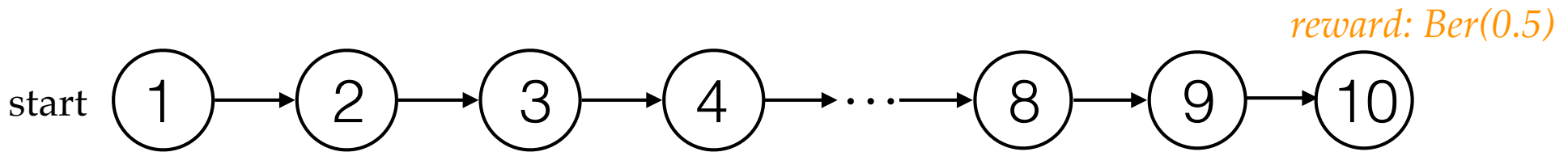
Bellman completeness (closure)

$$\mathcal{T} f \in \mathcal{F}, \forall f \in \mathcal{F}$$

*introduced by Szepesvari
& Munos [2005]

- whatever f_{k-1} is used, regression is always well-specified
- Implies realizability for finite class (why?)
- For piecewise const F , completeness = bisimulation (hw)
- Not necessarily converge, but will get close to a good solution (under additional data assumptions)

How completeness fixes the issue



Function
class

0.5 × 0.5
~~0.528~~ 0.628 0.1
 0.628

Tf_{t-1}

- More generally: issue goes away if the regression problem

$$\left\{ \left((s, a), \left(r + \gamma \max_{a' \in \mathcal{A}} f_{t-1}(s', a') \right) \right) \right\}$$

is realizable with F , for any $f_{t-1} \in F$

~~$Hf \in F$~~
 $Tf \in F$

- In **finite-horizon** setting: the richer function class you use at a lower level, the **more difficult** to satisfy realizability at higher level
- In **discounted** setting: F closed under Bellman update—adding functions can **hurt** representation

Alternative approach

$$R(s,c) + \gamma \mathbb{E}_{s'|s,a} [\max_{a'} f(s',a')]$$

- FQI is an **iterative** alg in its nature
 - not optimizing a **fixed objective function!**
 - objective changes as current f changes
- Alternative: minimize $\|f - \mathcal{T}f\|$ over $f \in F$
 - Is it equivalent to minimizing:

$$\min_{\theta \in \mathcal{D}} L(\theta)$$

~~$$(f - Q^*)^2$$~~

$$f = Q^* \Leftrightarrow f = \mathcal{T}f.$$

$$(s,a) \sim \mu.$$

$$r, s'$$

$$\mathbb{E}_{\substack{(s,a) \sim \mu \\ r \sim R(s,a) \\ s' \sim P(s,a)}} \left[\left(f(s,a) - (r + \gamma \max_{a'} f(s',a')) \right)^2 \right]$$

(omitted in the rest of slides)

$$\mathbb{E}[\cdot | s,a] = \mathcal{T}f.$$

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

Bellman error minimization

$$\mathbb{E}_{(s,a) \sim \mu} \left[\left(f(s,a) - (r + \gamma \max_{a'} f(s', a')) \right)^2 \right]$$

$Q^*(s,a) - (r + \gamma \dots)$

$$= \mathbb{E}_{(s,a) \sim \mu} \left[(f(s,a) - (\mathcal{T}f)(s,a))^2 \right] + \mathbb{E}_{(s,a) \sim \mu} \left[\left((\mathcal{T}f)(s,a) - (r + \gamma \max_{a'} f(s', a')) \right)^2 \right]$$

This part is what we want:
 $\|f - \mathcal{T}f\|$, with a weighted
 2-norm defined w/ ν

This part is annoying!

- Prefer “flat” f
- Q^* is not necessarily flat!
- 0 for deterministic transitions. Issue is only serious when env highly stochastic

Unbiased estimate
 “double sampling”

Baird '95

(s, a, r_A, s'_A)
 (s, a, r_B, s'_B)

Workaround #1

- For $(s, a) \sim \mu$, if we can obtain **2** i.i.d. copies of (r, s') (copy A & B):

$$\left(f(s,a) - \left(r_A + \gamma \max_{a' \in \mathcal{A}} f(s'_A, a') \right) \right) \left(f(s,a) - \left(r_B + \gamma \max_{a' \in \mathcal{A}} f(s'_B, a') \right) \right)$$

- Only doable in simulators w/ resets...

Bellman error minimization

if $Tf \in \mathcal{G}$ \rightarrow $\min_{g \in \mathcal{G}} \mathbb{E} [g(s,a) - (r + \gamma \max_{a'} f(s', a'))^2]$

$$\mathbb{E}_{(s,a) \sim \mu} \left[\left(f(s,a) - (r + \gamma \max_{a'} f(s', a')) \right)^2 \right]$$

$$= \mathbb{E}_{(s,a) \sim \mu} \left[(f(s,a) - (Tf)(s,a))^2 \right] + \mathbb{E}_{(s,a) \sim \mu} \left[\left((Tf)(s,a) - (r + \gamma \max_{a'} f(s', a')) \right)^2 \right]$$

This part is what we want:
 $\|f - Tf\|$, with a weighted
 2-norm defined w/ ν

This part is annoying!

- Prefer “flat” f
- Q^* is not necessarily flat!
- 0 for deterministic transitions. Issue is only serious when env highly stochastic

Workaround #2

- Estimate the 2nd part, and subtract it from LHS

- Antos et al'08: *argmin $f \in \mathcal{F}$ $L(f)$. where $L(f) :=$*

$$\mathbb{E}_{(s,a) \sim \mu} \left[\left(f(s,a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f(s', a') \right) \right)^2 \right] - \min_{g \in \mathcal{G}} \mathbb{E}_{(s,a) \sim \mu} \left[\left(g(s,a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f(s', a') \right) \right)^2 \right]$$

argmin over \mathcal{F}

$$Q^* \in \mathcal{F}$$

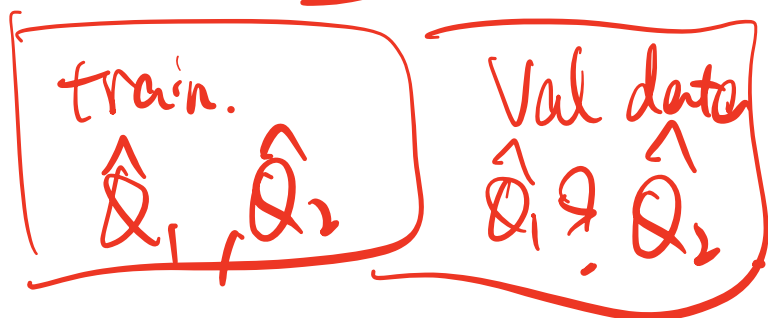
$$\forall f \in \mathcal{F}, \exists T f \in \mathcal{G}$$

$$\forall f \in \mathcal{F}, T f \in \mathcal{G}$$

Bellman error minimization

$$\arg \min_{f \in \mathcal{F}} \max_{g \in \mathcal{G}} \left(\mathbb{E}_{(s,a) \sim \mu} \left[\left(f(s,a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f(s',a') \right) \right)^2 - \left(g(s,a) - \left(r + \gamma \max_{a' \in \mathcal{A}} f(s',a') \right) \right)^2 \right] \right)$$

- Fix any f , the first squared error is constant; second square is a regression problem w/ Bayes optimal being Tf
- So, if \mathcal{G} is rich enough to contain Tf for all f , this works!
 - and w/ a consistent optimization objective, unlike FQI
- If \mathcal{G} is not rich enough, may under-estimate the Bellman error of some f (subtracting too much)
- FQI: When $\mathcal{G} = \mathcal{F}$, this is just Bellman completeness again!



$$\hat{Q}_1 - T\hat{Q}_1 ?$$

$$\hat{Q}_2 - T\hat{Q}_2 ?$$

$f_1 f_2.$

One last assumption: data

$\mathcal{T} \in \mathcal{F} \forall f \in \mathcal{F}.$

- Recall that data needs to be exploratory for batch RL
- What does it actually mean?
 - tabular: relatively uniform over state space
 - abstraction: relatively uniform over abstract state space
 - large/continuous state space: uniform? in what measure??

iid (s, a, r, s') .

Example.

$$r \sim R(s, a) \quad s' \sim P(\cdot | s, a)$$

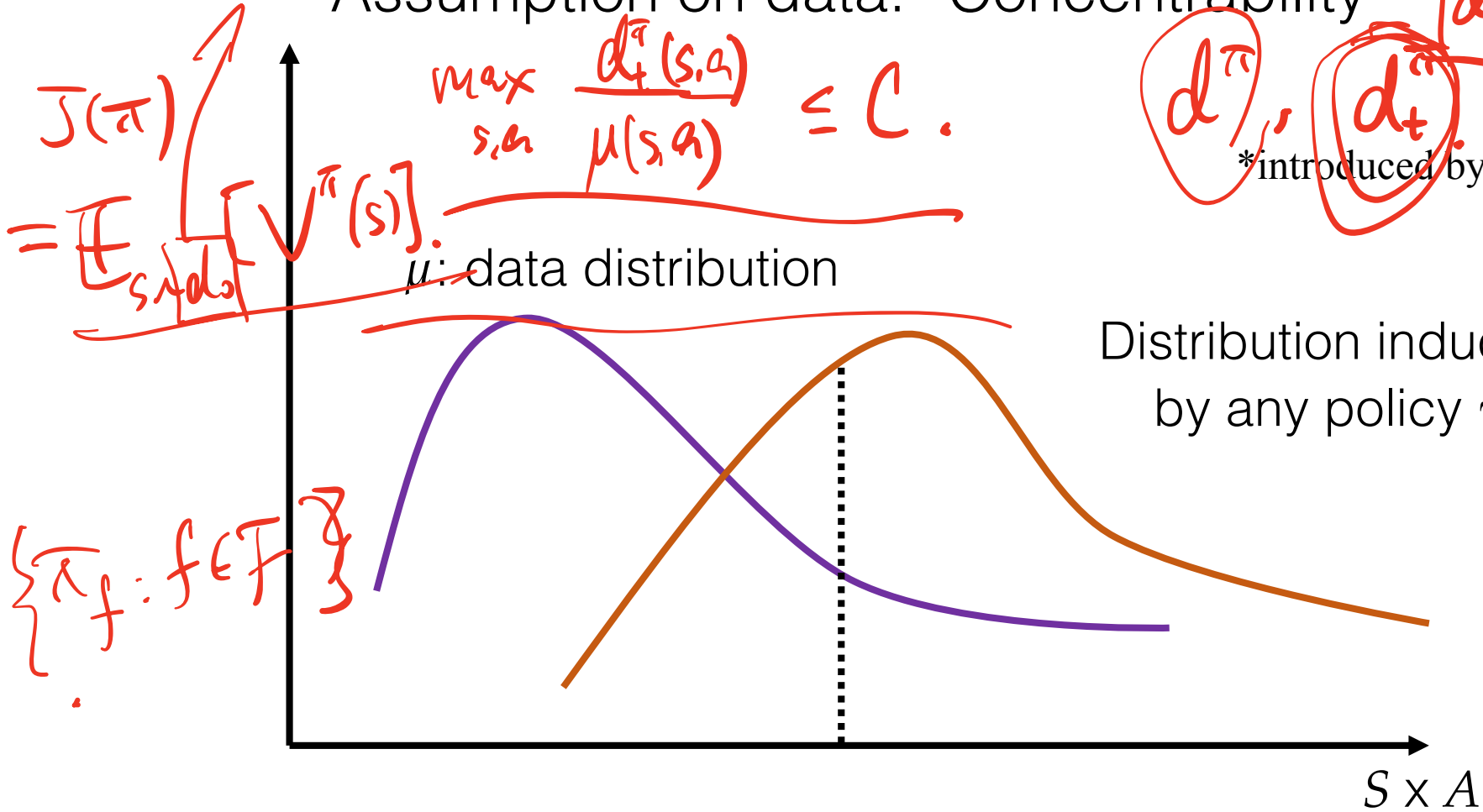
$$J(\pi^*) - J(\frac{A}{A})$$

Assumption on data: "Concentrability"

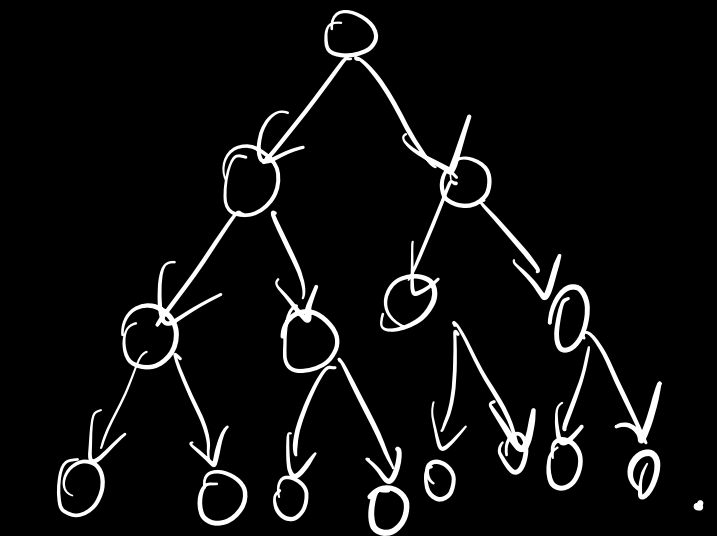
$$\max_{s,a} \frac{d_t^\pi(s,a)}{\mu(s,a)} \leq C.$$

d^π , d_t^π do.

*introduced by Munos'05

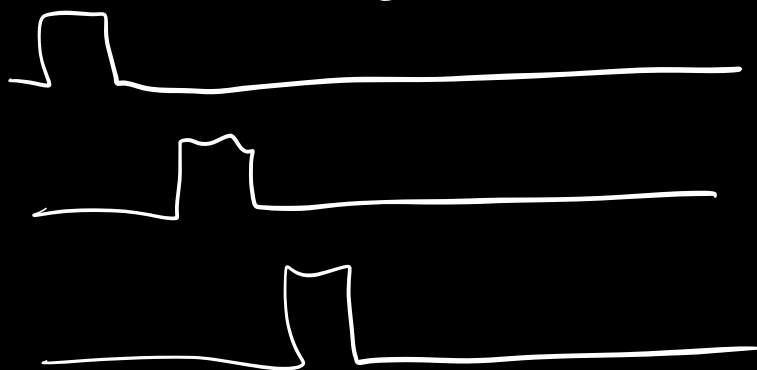


- Let C be a uniform upper bound on the density ratio
- Assumption: C is small (= allow polynomial dependence on C)
- Previous exponential lower bound is "explained away" by an exponentially large C



$$H \quad 2^H$$

$$\rightarrow t = H.$$

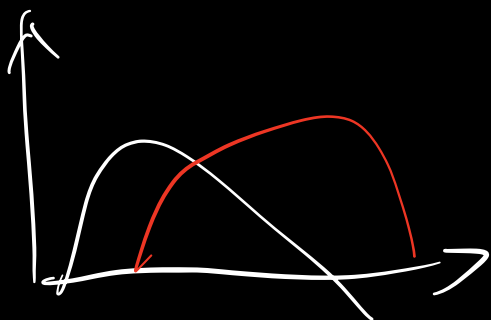


$$d_H^\pi$$

$$\max_{(s, \pi)} \frac{d_H^\pi(s)}{\mu_H(s)}$$

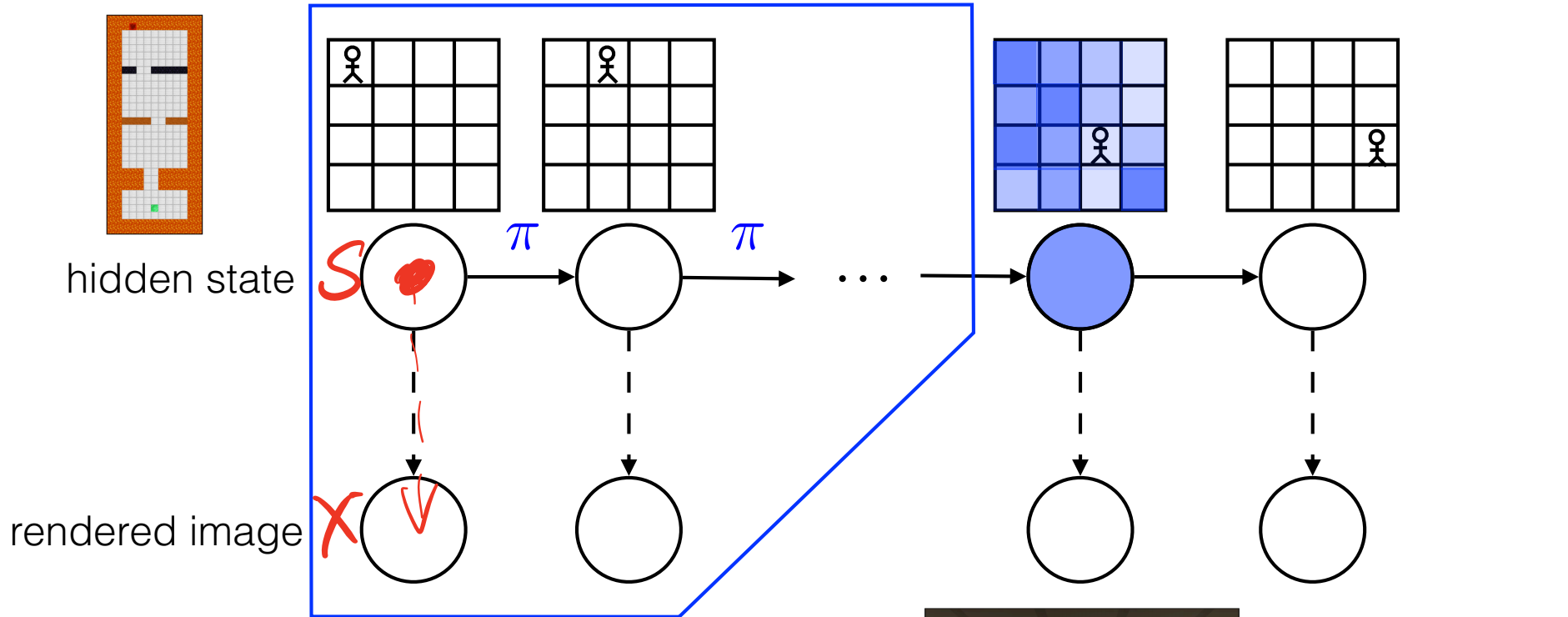
$$\leq \frac{2^H}{\dots}$$

$$= \max_{\pi} \left\| \frac{d_H^\pi}{\mu_H} \right\|_\infty.$$



$$\left\| \frac{d(\pi^*)}{\mu} \right\|_\infty$$

Concentrability: when is it small?



Connections to the assumptions
needed for efficient exploration
[Jiang et al'17]



Markovian high-dimensional observation

$$\frac{d^{\pi}(x)}{\mu^{\pi}(x)} \leq \frac{d^{\pi}(s) \cdot e[x|s]}{\mu(s) \cdot \rho(x|s)}$$

$\mu(s) \cdot e(x/s)$

Remainder of this part

Prove the $\text{poly}(H, \log|F|, C)$ result for FQI

Remainder of this part

Prove the $\text{poly}(H, \log|F|, C)$ result for FQI

FQE: $f_0 \equiv \vec{0}$.

$$\mathbb{E} \left[\begin{matrix} \uparrow \\ J f_{k-1} \\ |s, a \end{matrix} \right]$$

$$f_k = \underset{f \in \mathcal{F}}{\text{argmin}} \frac{1}{|D|} \sum_{(s, a, v, s')} (f(s, a) - v - \gamma \int f_{k-1}(s', \pi))^2$$

$$f_k \approx Q^\pi$$

$$J(\pi) = \mathbb{E}_{s \sim d_0} [Q^\pi(s, \pi)]$$

$$|J(\pi) - \hat{J}(\pi)|$$

$$\hat{J}(\pi) = \mathbb{E}_{s \sim d_0} [f_k(s, \pi)]$$

$$\hat{J} f \in \hat{\mathcal{F}}, \forall f \in \mathcal{F}$$

FQE: $f_0 \equiv \vec{0}$.

$$f_k = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{|D|} \sum_{(s, a, r, s')} \left(f(s, a) - r - \gamma \mathcal{T}^\pi f_{k-1}(s', \pi) \right)^2$$

(1) Want to show $f_k \approx \mathcal{T}^\pi f_{k-1}$.

$$\|f_k - \mathcal{T}^\pi f_{k-1}\|_2 \leq \text{small}$$

$$\Rightarrow \mathbb{E}_\mu \left[(f_k - \mathcal{T}^\pi f_{k-1})^2 \right] \leq \underline{\text{small}} \quad \forall k.$$

$$= \mathbb{E}_\mu \left[\left(\underbrace{f_k(s, a)}_{\mathcal{T}^\pi f_{k-1}(s, a)} - r - \gamma \underbrace{f_{k-1}(s', \pi)}_{f_{k-1}(s', \pi)} \right)^2 \right]$$

$$= \mathbb{E}_\mu \left[\left(\mathcal{T}^\pi f_{k-1}(s, a) - r - \gamma f_{k-1}(s', \pi) \right)^2 \right].$$

$$L_\mu(\cdot; \underline{f_{k-1}}) = \mathbb{E}_\mu \left[\left((\cdot) - r - \gamma \underline{f_{k-1}(s', \pi)} \right)^2 \right].$$

$$L_D(\cdot; f_{k-1}) = \frac{1}{|D|} \sum_{(s, a, r, s')} \left(\cdot - r - \gamma f_{k-1}(s', \pi) \right)^2.$$

fix any $\underline{f}, \underline{f}'$, $v_{\max} \geq 1 - \delta$.

$$|\underline{L}_{\mu}(\underline{f}; \underline{f}') - \underline{L}_D(\underline{f}; \underline{f}')| \leq O(v_{\max} \sqrt{\frac{1}{n} \ln \frac{1}{\delta}}).$$

$$f_k = \operatorname{argmin}_{f \in \mathcal{F}} L_D(f; f_{k-1}).$$

$$\mathcal{T}f_{k-1} = \operatorname{argmin}_{f \in \hat{\mathcal{F}}} L_{\mu}(f; f_{k-1}).$$

$$\mathbb{E}[(\underline{f}_k - \mathcal{T}f_{k-1})^2]$$

$$\leq \underline{L}_{\mu}(\underline{f}_k; f_{k-1}) -$$

$$L_D(\underline{f}_k; f_{k-1}) +$$

$$\underline{L}_{\mu}(\mathcal{T}f_{k-1}; f_{k-1}).$$

$$L_D(\mathcal{T}f_{k-1}; f_{k-1})$$

$$\leq 2 \cdot O(v_{\max} \sqrt{\frac{1}{n} \ln \frac{1}{\delta}}).$$

$$\mathbb{E}_{(s, a) \sim \mu} \left[\left(f_k(s, a) - (\mathcal{T}f_{k-1})(s, a) \right)^2 \right] \leq$$

$$\mathbb{E}_\mu[(\cdot)^2] = \|\cdot\|_{2,\mu}^2 \quad \|\cdot\|_{L^2(\mu)}$$

$$\left(\mathbb{E}_\mu[(\cdot)^p]\right)^{1/p} =: \|\cdot\|_{p,\mu}$$

$$\rightarrow \|\cdot\|_1 \not\approx \|\cdot\|_2 \not\approx \|\cdot\|_\infty$$

$$\|\cdot\|_1, \mu \leq \|\cdot\|_{2,\mu} \leq \|\cdot\|_{\infty,\mu}$$

if fully supp.

if $v \in \mathbb{R}^d$, chose $\mu = \text{unif}([d])$
 " $\{1, 2, \dots, d\}$.

$$\|v\|_1 = |v_1| + |v_2| + \dots + |v_d|$$

$$\begin{aligned} \|v\|_{1,\mu} &= \frac{1}{d} (|v_1| + |v_2| + \dots + |v_d|) \\ &= \frac{1}{d} \sum_i |v_i| \end{aligned}$$

$$\begin{aligned} \|v\|_{p,\mu} \leq \varepsilon & \quad \|v\|_{p,\mu} \leq \left(\frac{\mu}{\mu} \|v\|_\infty\right)^{1/p} \varepsilon \\ &= \left(\sum_{i=1}^d \mu(i) |v(i)|^p\right)^{1/p} \end{aligned}$$

$$|\hat{J}(\underline{\pi}) - J(\underline{\pi})|$$

$$= \left| \mathbb{E}_{s \sim d_0} [f_k(s, \underline{\pi})] - \mathbb{E}_{s \sim d_0} [Q^\pi(s, \underline{\pi})] \right|$$

$$= \left| \mathbb{E}_{\substack{s \sim d_0 \\ a \sim \pi}} [f_k(s, a)] - \left(\mathbb{E}_{\substack{a \sim \pi \\ r, s'}} [\gamma + \delta Q^\pi(s', \pi)] \right) \right|$$

$$+ \mathbb{E}_{\substack{s \sim d_0 \\ a \sim \pi \\ r, s'}} [\gamma + \delta f_{k-1}(s', \pi)]$$

$(s, a) \cdot I \leq \mathbb{E}_{d_0 \times \pi} [|f_k - T^\pi f_{k-1}|]$
 $\leq \|f_k - T^\pi f_{k-1}\|_{2, d_0 \times \pi}$
 $\leq \sqrt{\|d_0 \times \pi\|_\infty} \cdot \|f_k - T^\pi f_{k-1}\|_{2, \mu}$

$\| \mu \|_\infty \leq C$
 $\forall t$

$$(II) - (IV) = \delta \mathbb{E}_{s' \sim d_2^\pi} [f_{k-1}(s', \pi) - Q^\pi(s', \pi)]$$

	Data	Function approximation
AVI	$\max_{\pi} \ d^{\pi} / d^D\ _{\infty} \leq C$	$\mathcal{T}f \in \mathcal{F}, \forall f \in \mathcal{F}$ [Munos & Szepesvari '08]
API		$\mathcal{T}^{\pi} f \in \mathcal{F}, \forall f \in \mathcal{F}, \pi \in \Pi$ [Antos et al '08]

- **Assumption** so far: data is **exploratory** (e.g., $\max_{\pi} \|d^{\pi} / \mu\|_{\infty} \leq C$)
- Challenge: real-world data often **lacks** exploration!
 - Data may not contain all **bad** behaviors
 - Alg may **over-estimate** their performance



How to understand a driving behavior is unsafe, if all data are safe?

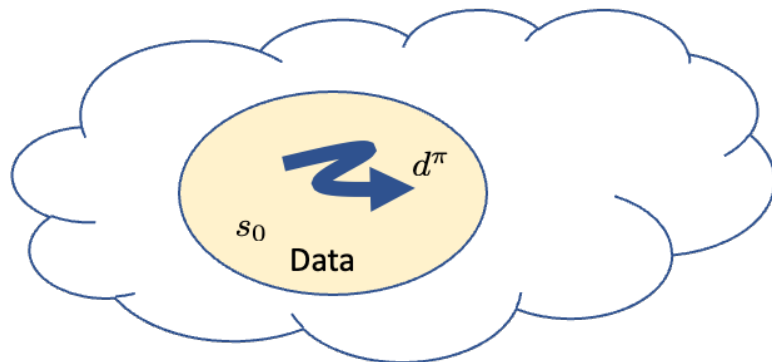
Data with **insufficient** coverage

- Policy optimization: $\arg \max_{\pi \in \Pi} J(\pi) := Q^\pi(s_0, \pi)$
 - Q^π : value function; s_0 : initial state; Π : policy class
- Considerations in **estimating** $\hat{J}(\pi)$?

$$\arg \max_{\pi \in \Pi} \hat{J}(\pi)$$

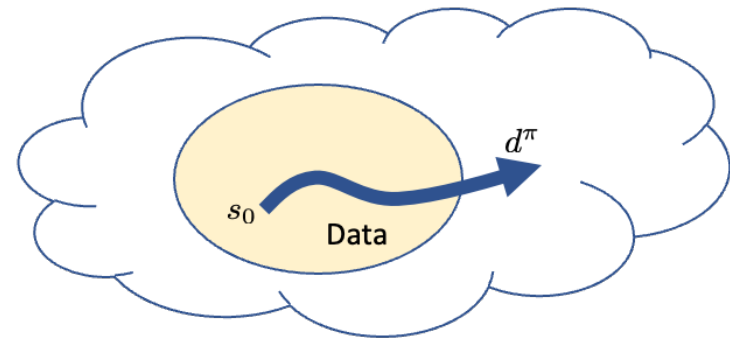
Pessimism in face of uncertainty

$$\hat{J}(\pi) \approx J(\pi)$$



Policy **covered** by data

$$\hat{J}(\pi) \leq J(\pi)$$



Policy **not covered** by data

Handle two cases simultaneously

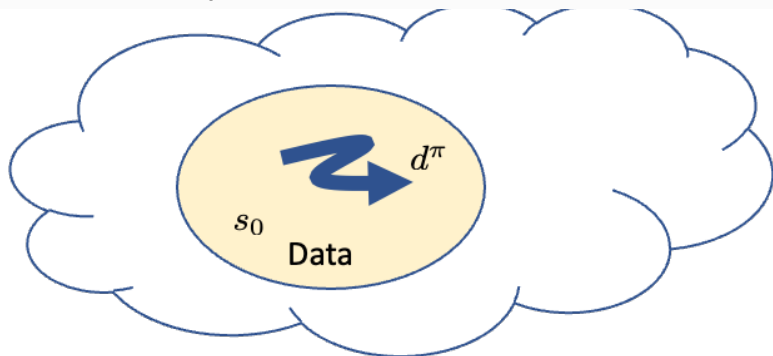
- Consider $\mathcal{F}_\epsilon^\pi := \{f \in \mathcal{F} : \|f - \mathcal{T}^\pi f\|_{2,\mu} \leq \epsilon\}$ “Confidence set”/“Version space”
 - small $\|f - \mathcal{T}^\pi f\|_{2,\mu}$ implies $f(s_0, \pi) \approx J(\pi) = Q^\pi(s_0, \pi)$ if μ covers d^π
 - can estimate $|f(s_0, \pi) - J(\pi)| \leq \frac{1}{1-\gamma} \|f - \mathcal{T}^\pi f\|_{2,d^\pi}$ (under “Bellman-completeness” $\mathcal{T}^\pi \mathcal{F} \in \mathcal{F}, \forall f \in \mathcal{F}$)
- Key observation:** Q^π is in the set ($Q^\pi - \mathcal{T}^\pi Q^\pi \equiv 0$)
- Pessimistic policy evaluation**

$$\hat{J}(\pi) := \min_{f \in \mathcal{F}_\epsilon^\pi} f(s_0, \pi) \leq Q^\pi(s_0, \pi) = J(\pi)$$

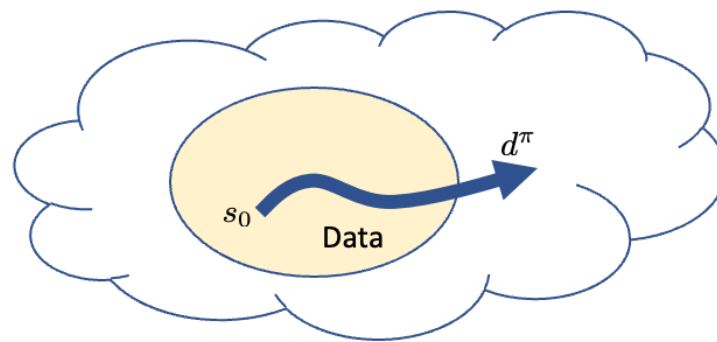


$$\hat{J}(\pi) \leq J(\pi)$$

All members of \mathcal{F}_ϵ^π have small $\|f - \mathcal{T}^\pi f\|_{2,\mu}$, so $\hat{J}(\pi) \approx J(\pi)$ for **covered** π



Policy **covered** by data



Policy **not covered** by data

	Data	Function approximation	
AVI	$\max_{\pi} \ d^{\pi} / d^D\ _{\infty} \leq C$	$\mathcal{T}f \in \mathcal{F}, \forall f \in \mathcal{F}$	[Munos & Szepesvari'08]
API		$\mathcal{T}^{\pi} f \in \mathcal{F}, \forall f \in \mathcal{F}, \pi \in \Pi$	[Antos et al '08]
Pessimism	$\ d^{\pi^*} / d^D\ _{\infty} \leq C$	$\mathcal{T}^{\pi} f \in \mathcal{F}, \forall f \in \mathcal{F}, \pi \in \Pi$	[Xie et al '21]

- Guarantee: $\hat{\pi} = \arg \min_{\pi \in \Pi} \hat{J}(\pi)$ competes with any **covered** policy $\pi_{\text{ref}} \in \Pi$
 - $J(\hat{\pi}) \geq \hat{J}(\hat{\pi}) \geq \hat{J}(\pi_{\text{ref}}) \approx J(\pi_{\text{ref}})$
 - **Near-optimality** follows if π^* is **covered**
- Alternative: **pointwise** pessimism (construct $\hat{Q}^{\pi}(s, a) \leq Q^{\pi}(s, a) \quad \forall s, a$)
 - Insert negative bonus in Bellman backup [Jin et al'21]
 - Density estimation + pessimistic in low-density area [Liu et al'20]