

Model-based approaches for continuous state

estimating a model & then using to compute a V or policy
 often considered more sample efficient than model-free estimates
 likely to be true for model-free algorithms that only use data
 tuples (s, a, r, s') once like Q-learning
 but when use experience replay or multiple iterations through
 the data, less clear
 as in value function approximation, many ways to approximate a
 model

- Approximate Linear Models & Connection to LSP!
- assume set of indep features for representing transition & reward models
 $\Phi(s) = [\phi_1(s) \dots \phi_K(s)]^T$ k features
 $\Phi(s'|s)$ for a fixed π
 1st consider uncontrolled / fixed policy case
 $P(s'|s, a) \rightarrow P(s'|s)$ shorthand because only one action per s
- define $\Phi(s'|s) = \sum_{s' \sim P(s'|s)} [\phi_1(s') \dots \phi_K(s')]$
 want to compute a $k \times k$ matrix P_Φ (like w for value func approx) that
 1) predicts expected next feature vectors
 $P_\Phi^T \Phi(s) \approx E_{s' \sim P(s'|s)} \{ \Phi(s'|s) \}$ $k \times 1$
 2) minimizes the expected feature-prediction error
 $P_\Phi = \underset{P_K}{\operatorname{argmin}} \sum_s \| P_K^T \Phi(s) - E \{ \Phi(s'|s) \} \|_2^2$
 solving the optimization in (2)
 compute the expected next state directly as P_Φ \downarrow true transition model
 P_Φ is $n \times k$ matrix
 its row is expected value of features on next step after starting in s_i
 Φ is $n \times k$ matrix of feature values for each state
 its row of $P_\Phi \Phi$ is P_Φ 's prediction of next feature values for state i
 $\Phi P_\Phi \approx P_\Phi$
 least squares soln : $\Phi^T P_\Phi = \Phi^T P_\Phi$
 $P_\Phi = (\Phi^T \Phi)^{-1} \Phi^T P_\Phi$
 yields a next feature vec $\hat{P}_\Phi = \Phi \hat{r}_\Phi$
 predict reward using some features, project (using least squares)
 reward into space of features
 $r_\Phi = (\Phi^T \Phi)^{-1} \Phi^T R$
 $\hat{R} = \Phi \hat{r}_\Phi$

now prove linear fixed point soln for \mathbb{E} (resulting w)
identical to soln for value get using $P_{\mathbb{E}}^T r_{\mathbb{E}}$

- uncontrolled setting
let x be any k -length vector (represents a state)

Bellman eqn

$$V[x] = r_{\mathbb{E}}^T x + \gamma V[P_{\mathbb{E}}^T x]$$

$$= \sum_{i=0}^{\infty} \gamma^i r_{\mathbb{E}}^T (P_{\mathbb{E}}^T)^i x$$

can rewrite in terms of original state space as

$$V = \underbrace{\mathbb{E} \sum_{i=0}^{\infty} \gamma^i P_{\mathbb{E}}^T}_{\text{some linear comb of } \mathbb{E}'s \text{ columns}} r_{\mathbb{E}}$$

implies can express $V = \mathbb{E}w'$ for some w vector w

$$V = \hat{R} + \gamma \hat{P}_{\mathbb{E}} w'$$

$$\mathbb{E}w' = \hat{R} + \gamma \hat{P}_{\mathbb{E}} w'$$

$$\mathbb{E}w' = \mathbb{E}r_{\mathbb{E}} + \gamma \mathbb{E}P_{\mathbb{E}} w \quad \text{sub in expr for } \hat{P}_{\mathbb{E}} + \hat{R}$$

$$w' = (I - \gamma P_{\mathbb{E}})^{-1} r_{\mathbb{E}}$$

\nwarrow identity matrix

$\underbrace{\quad \quad \quad}_{\text{well defined if } P_{\mathbb{E}} \text{ has a spectral radius } < 1/\gamma}$
(if $> 1/\gamma$ implies value of some states is unbounded)

thm: for any Markov reward process and a set of features \mathbb{E}
the linear model soln & linear fixed point soln are identical

proof:

soln for linear model is

$$w' = (I - \gamma P_{\mathbb{E}})^{-1} r_{\mathbb{E}}$$

$$= (I - \gamma (\mathbb{E}^T \mathbb{E})^{-1} \mathbb{E}^T P_{\mathbb{E}})^{-1} (\mathbb{E}^T \mathbb{E})^{-1} \mathbb{E}^T R$$

$$= w_{\mathbb{E}} \quad \text{for fixed point} \quad \square$$

\Rightarrow for a given set of features \mathbb{E} , get exact same soln if compute approx linear model & then use to compute value function,
as directly computing approx linear fixed point value func

- controlled case: LSP!

$Q^n(s, a)$, define basis functions over $s \times a$ space

$$\hat{Q} = \sum_{i=1}^k w_i \phi_i(s, a)$$

policy eval: use LSTDQ to compute \hat{Q} of a fixed π
a MDP with a fixed π is equal to a MRP whose state space is SAA

LSTDQ = LSTD executed on this equiv MRP

by thm above, LSTDQ soln = approx linear models w/smallest L_2 error