

Recent Advances of Statistical Reinforcement Learning

Part 2

Gergely Neu (Universitat Pompeu Fabra)

Sattar Vakili (MediaTek Research)

Tutorial, UAI 2024

Part 2

1. Introduction to structural complexity
2. Linear Function Approximation
3. Non-linear Function Approximation

Tabular Setting

Recall results for the tabular setting:

- Q-learning with UCB: [Jin et al., 2018]

$$\text{Regret}(T) = O(\overline{H^3 SAT})$$

- Sample complexity:

$$\tilde{O}\left(\frac{\text{poly}(H)SA}{2}\right)$$

Tabular Setting

Recall results for the tabular setting:

- Q-learning with UCB: [Jin et al., 2018]

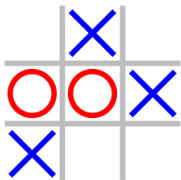
$$\text{Regret}(T) = O(\sqrt{H^3 S A T})$$

- Sample complexity:

$$\tilde{O}\left(\frac{\text{poly}(H) S A}{\epsilon}\right)$$

These are only meaningful if $T \geq S$ or $\epsilon \geq 1/\sqrt{S}$!

Why “Tabular”?



- Small size of state-action space
- $Q(s, a)$ can be represented as a [table](#)

Why Function Approximation?

Number of states S is **enormous** in real-world problems!



- Game of Go: 10^{170} states
- Atari: 10^{100} states
- Physical systems:
continuum of states

Why Function Approximation?

Two types of challenges:

- | **Computational:** Q and V cannot even be stored in memory, and Bellman equations are intractable to solve even if P and r were known
- | **Statistical:** Most states are not visited even once! How could we expect to learn about P or r like that?

Why Function Approximation?

Two types of challenges:

| **Computational:** Q and V cannot even be stored in memory, and Bellman equations are intractable to solve even if P and r were known

| **Statistical:** Most states are not visited even once! How could we expect to learn about P or r like that?

We need to find a way to **generalize** knowledge from visited states to unvisited states by leveraging **structure**

RL with Function Approximation

- | Approximate value function $Q(s, a)$ (or policy) in a class F .
- | Hope that F captures the MDP structure appropriately and leverage the information in structure of F to learn faster if possible.
- | Typical function classes: Linear, Kernel-based, NN-based

Tabular Linear Nonlinear

Setting

- | Generative oracle, Offline, Online
- | Episodic, Infinite horizon (discounted)
- | Model-based, Model-free

In this part we focus on:

Tabular Linear Nonlinear

Setting

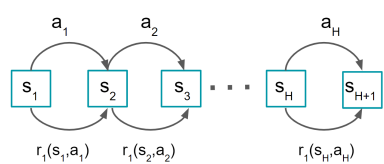
- | Generative oracle, Offline, Online
- | Episodic, Infinite horizon (discounted)
- | Model-based, Model-free

In this part we focus on:

Tabular Linear Nonlinear

Setting

For a clear and sharp presentation we focus on **episodic MDPs**

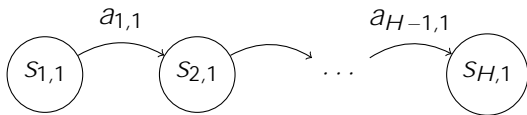


For simplicity, we assume r is known and deterministic

We focus on the structural complexity of $P(s/s, a)$

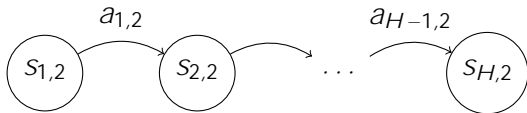
Episodic MDP

Episode 1:

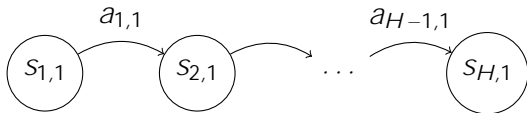


Episodic MDP

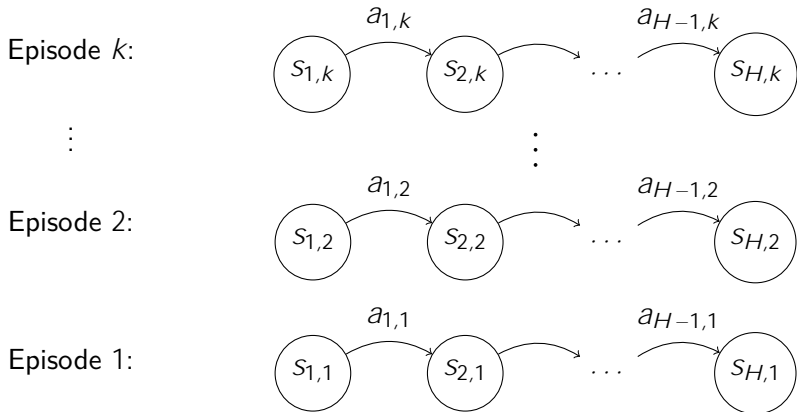
Episode 2:



Episode 1:



Episodic MDP



Part 2

1. Introduction to structural complexity
2. Linear Function Approximation
3. Non-linear Function Approximation

Linear Function Approximation

IDEA: approximate the Q -functions as linear functions of a given d -dimensional feature map $\phi : S \times A \rightarrow \mathbb{R}^d$.

Linear Function Approximation

IDEA: approximate the Q -functions as linear functions of a given d -dimensional feature map $\phi : S \times A \rightarrow \mathbb{R}^d$.

Let $\Phi \in \mathbb{R}^{(S \times A) \times d}$ be the “matrix” of stacked feature vectors $[\phi(s_1, a_1) \dots \phi(s_N, a_N)]$ (where $N = |S \times A|$).

We need to find a parameter vector w such that $Q \approx \Phi w$.
(Meaning that $Q(s, a) \approx \phi(s, a) \cdot w$.)

Linear Function Approximation

IDEA: approximate the Q -functions as linear functions of a given d -dimensional feature map $\phi : S \times A \rightarrow \mathbb{R}^d$.

Let $\Phi \in \mathbb{R}^{(S \times A) \times d}$ be the “matrix” of stacked feature vectors $[\phi(s_1, a_1) \dots \phi(s_N, a_N)]$ (where $N = |S \times A|$).

We need to find a parameter vector w such that $Q \approx \Phi w$.
(Meaning that $Q(s, a) \approx \phi(s, a) \cdot w$.)

QUESTION: When can we learn w efficiently?

Linear Function Approximation

Various conditions on the feature map have been studied:

- **Linear Q** : there exists a such that $Q =$.
- **Linear Q** : for every policy , there exists a such that $Q =$.
- **Closure under Bellman operator**: for any $Q_\theta =$, $TQ_\theta \in \text{span}()$.
- **Linear MDP**: The transition and reward functions are linear in the features. This implies all of the above conditions.

A number of more refined conditions have been also studied, such as assuming linearity of V in some feature map, or other types of factorized transition models. We refer to [Du et al. \[2021\]](#), [Jin et al. \[2021\]](#) for more details on such extensions.

What can we hope for?

WANT: find an ϵ -optimal policy with a computational and sample complexity polynomial in d , $1/\epsilon$ and H , independently of S and A .

What can we hope for?

WANT: Find an ϵ -optimal policy with a computational and sample complexity polynomial in d , $1/\epsilon$ and H , independently of S and A .

This is **impossible** when only requiring linear realizability!
[Weisz et al., 2021]

What can we hope for?

WANT: Find an ϵ -optimal policy with a computational and sample complexity polynomial in $d, 1/\epsilon$ and H , independently of S and A .

This is **impossible** when only requiring linear Q -realizability!
[Weisz et al., 2021]

Polynomial sample complexity is possible when relaxing the condition to linear Q -realizability, but no practical algorithms are known [Weisz et al., 2022, 2023].

What can we hope for?

WANT: Find an ϵ -optimal policy with a computational and sample complexity polynomial in $1/\epsilon$ and H , independently of S and A .

This is **impossible** when only requiring linear Q -realizability!
[Weisz et al., 2021]

Polynomial sample complexity is possible when relaxing the condition to linear Q -realizability, but no practical algorithms are known [Weisz et al., 2022, 2023].

Situation is similar when only assuming closure under Bellman operator / Bellman completeness [Zanette et al., 2020b, Du et al., 2021].

What can we hope for?

WANT: Find an ϵ -optimal policy with a computational and sample complexity polynomial in $1/\epsilon$ and H , independently of S and A .

This is **impossible** when only requiring linear Q -realizability!
[Weisz et al., 2021]

Polynomial sample complexity is possible when relaxing the condition to linear Q -realizability, but no practical algorithms are known [Weisz et al., 2022, 2023].

Situation is similar when only assuming closure under Bellman operator / Bellman completeness [Zanette et al., 2020b, Du et al., 2021].

Linear MDP condition enables both statistical and computational efficiency!!! [Jin et al., 2023]

Linear MDPs

Linear transition function:

$$P_h(j|s; a) = \sum_i h(s; a)_i h(j)_i;$$

where $h(\cdot) = [h^1(\cdot); \dots; h^d(\cdot)]$ is a d -dimensional signed measure.

Linear rewards: $r_h(s; a) = \sum_i h(s; a)_i w_i$.

Linear MDPs

Linear transition function:

$$P_h(j|s; a) = \int h(s'; a) \mu_h(s, j, a);$$

where $\mu_h(s, j, a) = [\mu_h^1(s, j, a); \dots; \mu_h^d(s, j, a)]$ is a d -dimensional signed measure.

Linear rewards: $r_h(s; a) = \int h(s'; a) \#_h ds'$.

In matrix notation:

Transition operator $P_h \in \mathbb{R}^{(S \times A) \times S}$ can be written as

$P_h = M_h$ for some "matrix" $M_h \in \mathbb{R}^{S \times d}$.

Reward function can be written as $r_h = \#_h$ for some $\#_h \in \mathbb{R}^d$.

Tabular MDPs are linear

Tabular setting is a special case with
dimension $d = SA$:

Let $\phi_h(s; a) = e_{(s;a)}$ be the
canonical basis in \mathbb{R}^d

$$P_h(j; s; a) = e_{s;a}^T \phi_h(j)$$

A magical property of linear MDPs

In a linear MDP, the Q-functions of all policies are linear in:

$$\begin{aligned} Q_h &= r_h + P_h V_{h+1} = \#_h + M_h V_{h+1} \\ &= \#_h + M_h V_{h+1} = \tilde{Q}_h; \end{aligned}$$

with $\tilde{Q}_h = \#_h + M_h V_{h+1}$.

A magical property of linear MDPs

In a linear MDP, the Q-functions of all policies are linear in:

$$\begin{aligned} Q_h &= r_h + P_h V_{h+1} = \#_h + M_h V_{h+1} \\ &= \#_h + M_h V_{h+1} = \tilde{Q}_h; \end{aligned}$$

with $\tilde{Q}_h = \#_h + M_h V_{h+1}$.

This implies linear Q^* -realizability, linear Q -realizability, Bellman completeness, and many more useful properties for analysis! E.g., note that for any function $u \in \mathbb{R}^S$, $P_h u = M_h u$ is linear in u .

A magical property of linear MDPs

In a linear MDP, the Q-functions of all policies are linear in:

$$\begin{aligned} Q_h &= r_h + P_h V_{h+1} = \#_h + M_h V_{h+1} \\ &= \#_h + M_h V_{h+1} = \tilde{Q}_h; \end{aligned}$$

with $\tilde{Q}_h = \#_h + M_h V_{h+1}$.

This implies linear Q[?]-realizability, linear Q-realizability, Bellman completeness, and many more useful properties for analysis! E.g., note that for any function $u \in \mathbb{R}^S$, $P_h u = M_h u$ is linear in u .

The structure of linear MDPs allows us to import tools from linear bandit literature [Abbasi-Yadkori et al., 2011, Lattimore and Szepesvári, 2020].

Optimistic approximate dynamic programming

IDEA: Combine the techniques for tabular MDPs with exploration bonuses borrowed from the linear bandit literature!

Optimistic approximate dynamic programming

IDEA: Combine the techniques for tabular MDPs with exploration bonuses borrowed from the linear bandit literature!

UCB-VI [Azar et al., 2017]:

Backtrack $h = H; H - 1; \dots; 1$: run optimistic value iteration

$$Q_h = r_h + \underbrace{\frac{b_h}{|Z_h|}}_{\text{model estimate}} V_{h+1} + \underbrace{\frac{b_h}{|Z_h|}}_{\text{exploration bonus}}$$

and set $V_h(s) = \max_a Q_h(s; a)$ for all $s; a$.

Forward $h = 1; 2; \dots; H$: take actions according to greedy policy

$$\pi_h(s) = \arg \max_a Q_h(s; a):$$

Optimistic approximate dynamic programming

IDEA: Combine the techniques for tabular MDPs with exploration bonuses borrowed from the linear bandit literature!

UCB-VI [Azar et al., 2017]:

Backtrack $h = H; H - 1; \dots; 1$: run optimistic value iteration

$$Q_h = r_h + \underbrace{\frac{b_h}{|Z_h|}}_{\text{model estimate}} V_{h+1} + \underbrace{\frac{b_h}{|Z_h|}}_{\text{exploration bonus}}$$

and set $V_h(s) = \max_a Q_h(s; a)$ for all $s; a$.

Forward $h = 1; 2; \dots; H$: take actions according to greedy policy

$$\pi_h(s) = \arg \max_a Q_h(s; a)$$

But how do we define P_h and b_h ?

Least Squares Value Iteration (LSVI)

Transition model P_h can be defined **implicitly** via least-squares:

- | Solve the regularized linear regression problem

$$\mathbf{w}_{h;k} = \arg \min_{\mathbf{w}} \sum_{t=1}^k (V_{h+1;k}(s_{h+1;t}) - h(s_{h;t}; \mathbf{a}_{h;t}); \mathbf{w})^2 + \lambda \|\mathbf{w}\|^2$$

- | That provides a prediction

$$[P_h^{\setminus} V_{h+1;k}](s; \mathbf{a}) = h(s; \mathbf{a}); \mathbf{w}_{h;k}$$

- | Also, an uncertainty quantification (variance)

$$\sigma_{h;k}^2(s; \mathbf{a}) = \frac{\mathbf{k}(s; \mathbf{a}) \mathbf{k}^2}{\lambda + \sum_{t=1}^k \chi_{h;t}^2}$$

$$\chi_{h;t}^2 = \sum_{s, \mathbf{a}} \mathbb{1}_{\{s_{h;t} = s, \mathbf{a}_{h;t} = \mathbf{a}\}} (s_{h;t}; \mathbf{a}_{h;t}) (s_{h;t}; \mathbf{a}_{h;t})$$

LSVI-UCB

The prediction and variance give us an upper confidence bound on Q

$$Q_{h;k}(s; a) = r_h(s; a) + [P_h V_{h+1}](s; a) + \beta_{h;k}(s; a)$$

This is then used to compute an UCB on Q as

$$V_{h;k}(s) = \max_a Q_{h;k}(s; a)$$

Performance guarantees

Theorem [Jin et al., 2023] The regret of LSVI-UCB satisfies
$$\text{Regret}(K) = \mathcal{O}(H^2 \sqrt{d^3 K}).$$

This implies a sample complexity guarantee of $\mathcal{O}\left(\frac{\text{poly}(H)d^3}{\epsilon}\right)$.

Performance guarantees

Theorem [Jin et al., 2023] The regret of LSVI-UCB satisfies

$$\text{Regret}(K) = \mathcal{O}\left(H^2 \sqrt{d^3 K}\right).$$

This implies a sample complexity guarantee $\mathcal{O}\left(\frac{\text{poly}(H)d^3}{\epsilon}\right)$.

Proof ideas:

Prove confidence bounds

$$\mathbb{P}\left[\left|P_{h,k} \setminus V_{h+1;k}(s; a) - [P_{h,k} V_{h+1;k}](s; a)\right| \leq \beta_{h,k}(s; a)\right]$$

Using standard techniques (e.g., Azar et al., 2017), show

$$\text{Regret}(K) \leq \sum_{h,k} \mathbb{P}\left[\sum_{h,k} \beta_{h,k}(s_{h,k}; a_{h,k})\right]$$

Use elliptical potential lemma (e.g., Abbasi-Yadkori et al., 2011) to show

$$\sum_{h,k} \mathbb{P}\left[\sum_{h,k} \beta_{h,k}(s_{h,k}; a_{h,k})\right] \leq H \sqrt{Kd \log(K)}$$

Proving the confidence bounds

By standard results on least-squares estimators (e.g., [Abbasi-Yadkori et al., 2011](#)), one can prove the following confidence bound for any $\mathbf{u} \in \mathbb{R}^S$ that holds with probability at least $1 - \delta$:

$$|\hat{P}_h \mathbf{u}(s; \mathbf{a}) - P_h \mathbf{u}(s; \mathbf{a})| \leq \beta_{h,k}(s; \mathbf{a})$$

for some

$$\beta_{h,k}(s; \mathbf{a}) = \frac{1}{2} k M_h \|\mathbf{u}\| + H^q \frac{1}{d \log(\frac{K}{\delta})}$$

Proving the confidence bounds

By standard results on least-squares estimators (e.g., [Abbasi-Yadkori et al., 2011](#)), one can prove the following confidence bound for any $\mathbf{u} \in \mathbb{R}^S$ that holds with probability at least $1 - \delta$:

$$|[\hat{\mathbf{P}}_h \mathbf{u}](s; \mathbf{a}) - [\mathbf{P}_h \mathbf{u}](s; \mathbf{a})| \leq \beta_{h,k}(s; \mathbf{a})$$

for some

$$\beta_{h,k}(s; \mathbf{a}) = \frac{1}{2} k M_h \|\mathbf{u}\| + H^q \sqrt{\frac{\log(\frac{K}{\delta})}{d}}$$

Challenge: $\mathbf{u} = \mathbf{V}_{h+1;k}$ is **not** fixed, but depends on all past data!

Proving the confidence bounds

By standard results on least-squares estimators (e.g., [Abbasi-Yadkori et al., 2011](#)), one can prove the following confidence bound for any $\mathbf{u} \in \mathbb{R}^S$ that holds with probability at least $1 - \delta$:

$$|\hat{P}_{h,\mathbf{u}}(\mathbf{s}; \mathbf{a}) - P_{h,\mathbf{u}}(\mathbf{s}; \mathbf{a})| \leq \sqrt{\frac{1}{\gamma} \sum_{k=1}^{\gamma} \mathbf{u}^\top \mathbf{V}_{h+1;k}(\mathbf{s}; \mathbf{a}) \mathbf{u}}$$

for some

$$\frac{1}{\gamma} \sum_{k=1}^{\gamma} \mathbf{u}^\top \mathbf{V}_{h+1;k}(\mathbf{s}; \mathbf{a}) \mathbf{u} \leq \frac{1}{2} \kappa M_h \|\mathbf{u}\|^2 + H^q \frac{1}{d \log(\frac{\kappa}{\delta})}$$

Challenge: $\mathbf{u} = \mathbf{V}_{h+1;k}$ is **not** fixed, but depends on all past data!

Solution: Covering number argument

Covering Number Argument

- Notice that all value functions $V_{h,k}$ belong to the function class

$$V = \{V(s) = \min_{f \in H} \max_{a \in A} (s; a) + k(s; a) \mid g\}$$

- Idea:** cover the space of functions V such that we can rewrite

$$[P_h V_{h+1; k}](s; a) = [P_h V_{h+1; k}](s; a) + \sup_{u \in V} [P_h u](s; a) - [P_h u](s; a)$$

- How many functions u are required to cover V up to error?

$$N = \mathcal{O}(d^2)$$

Covering Number Argument

- Notice that all value functions $V_{h,k}$ belong to the function class

$$V = \{V(s) = \min_{f \in H} \max_{a \in A} (s; a) + k(s; a) \mid g\}$$

- Idea:** cover the space of functions V such that we can rewrite

$$[P_h V_{h+1; k}](s; a) = [P_h V_{h+1; k}](s; a) + \sup_{u \in V} [P_h u](s; a) - [P_h u](s; a)$$

- How many functions u are required to cover V up to error?

$$N = \mathcal{O}(d^2)$$

Covering Number Argument

- Notice that all value functions $V_{h,k}$ belong to the function class

$$V = \{V(s) = \min_{f \in H} \max_a \sum_{t=0}^{\infty} \gamma^t (f(s; a) - k(s; a)) \mid g\}$$

- Idea:** cover the space of functions V such that we can rewrite

$$[P_h V_{h+1; k}](s; a) - [P_h V_{h+1; k}](s; a) + \sup_{u \in V} [P_h u](s; a) - [P_h u](s; a) :$$

- How many functions u are required to cover V up to error?

$$N = \mathcal{O}(d^2)$$

Covering Number Argument

I We can now use a union-bound argument to show that

$$\sup_{u \in \mathcal{V}} |\mathbb{P}_h[u](s; a) - \mathbb{P}[u](s; a)| + \mathbb{P}(\|V_{h;k}\| \geq \epsilon)$$

holds with probability at least $1 - \epsilon$.

I Choosing $\epsilon = \frac{1}{2} \epsilon$, we get

$$\mathbb{P}(\|V_{h;k}\| \geq \frac{1}{2} \epsilon) \leq \frac{1}{2} \epsilon + H^q \frac{1}{\epsilon} \log\left(\frac{KN}{\epsilon}\right) = \mathcal{O}(Hd)$$

Alternative linear models

Linear MDP model factorizes $P = M$ with some known $M \in \mathbb{R}^{(S \times A) \times d}$ and some unknown $M \in \mathbb{R}^{d \times S}$.

Some alternative factorizations are:

Linear mixture MDPs [Zhou et al., 2021]: $P = M$ with some known $M \in \mathbb{R}^{(S \times A \times S) \times d}$ and unknown $M \in \mathbb{R}^{d \times S}$. Analysis is simpler but the model doesn't allow simple and explicit Q -function approximation and leads to impractical algorithms.

MatrixRL [Yang and Wang, 2020]: $P = M$ with some known $M \in \mathbb{R}^{(S \times A) \times m}$, another known $M \in \mathbb{R}^{n \times S}$, and an unknown $M \in \mathbb{R}^{m \times n}$. Can be shown to be a special case of linear mixture MDPs, and suffers from the same limitations.

Low-rank MDPs [Modi et al., 2024]: Same as linear MDPs except both M and M are unknown and belong to finite model class. Requires much more sophisticated techniques, but algorithms are kind of tractable.

Some References

Linear MDPs: [Jin et al. \[2020, 2023\]](#), [Yang and Wang \[2019, 2020\]](#), [Neu and Pike-Burke \[2020\]](#)

Linear Bellman complete models: [Zanette et al. \[2020a\]](#)

Linear mixture MDPs: [Yang and Wang \[2020\]](#), [Ayoub et al. \[2020\]](#), [Zhou et al. \[2021\]](#), [Moulin and Neu \[2023\]](#)

Other model classes with hidden finite-dimensional linear structure: [Du et al. \[2021\]](#), [Jin et al. \[2021\]](#)

Part 2

1. Introduction to structural complexity
2. Linear Function Approximation
3. Non-linear Function Approximation

Limitations of the Linear Setting

Directly reachable states:

$$S_{s;a} := \{s^0 \in S : P(s^0; a) > 0\}$$

$$U := \max_{(s;a) \in S \times A} |S_{s;a}|$$

Limitations of the Linear Setting

Directly reachable states:

$$S_{s;a} := \{s' \in S : P(s'|s; a) > 0\}$$

$$U := \max_{(s;a) \in S \times A} |S_{s;a}|$$

Theorem Lee and Oh [2024] For an MDP with a finite state space, the feature dimension is lower bounded by

$$d \geq \frac{|S|}{U} c$$

Limitations of the Linear Setting

High dimensional problems

Nonlinear problems

Kernel-Based Setting

Kernel-Based Setting

Kernel-based models are natural extensions of linear models to **in finite dimensional** feature maps

Kernel-Based Setting

Kernel-based models are natural extensions of linear models to **in finite dimensional** feature maps

Allow for versatile and powerful **non-linear function approximation**

Kernel-Based Setting

Kernel-based models are natural extensions of linear models to **infinite dimensional** feature maps

Allow for versatile and powerful **non-linear function approximation**

Lend themselves to analysis

Kernel-Based Setting

Kernel-based models are natural extensions of linear models to **infinite dimensional** feature maps

Allow for versatile and powerful **non-linear function approximation**

Lend themselves to analysis

Serve as an intermediate step towards analysis of NN-based models

Kernel-Based Setting

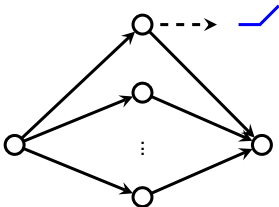
Kernel-based models are natural extensions of linear models to **in finite dimensional** feature maps

Allow for versatile and powerful **non-linear function approximation**

Lend themselves to analysis

Serve as an intermediate step towards analysis of NN-based models

Tabular! Linear! Kernel-Based NN-Based



Kernel-Based Setting

Function class:

$$\mathcal{F} = \left\{ f : \mathbb{R}^d \rightarrow \mathbb{R}; f(\cdot) = \sum_{m=1}^P w_m \phi_m(\cdot) \right\}$$

Kernel-Based Setting

Function class:

$$F = \{ f : \mathbb{R}^d \rightarrow \mathbb{R}; f(\cdot) = \sum_{m=1}^P w_m \phi_m(\cdot) \}$$

An extension of linear models to infinite dimensions in the feature space

Kernel-Based Setting

Function class:

$$F = \{ f : \mathbb{R}^d \rightarrow \mathbb{R}; f(\cdot) = \sum_{m=1}^P w_m \phi_m(\cdot) \}$$

An extension of linear models to infinite dimensions in the feature space

Nonlinear functions in \mathbb{R}^d

Mercer Theorem

A positive definite kernel $k : Z \times Z \rightarrow \mathbb{R}$

Mercer Theorem

A positive definite kernel $k: Z \times Z \rightarrow \mathbb{R}$

Theorem Any positive definite kernel can be written as

$$k(z; z^0) = \sum_{m=1}^{\infty} \lambda_m \phi_m(z) \phi_m(z^0)$$

The feature map $\phi_m(\cdot) = \sqrt{\lambda_m} \phi_m(\cdot)$ corresponding to

λ_m are referred to as eigenvalues

ϕ_m are referred to as eigenfunctions

Kernels

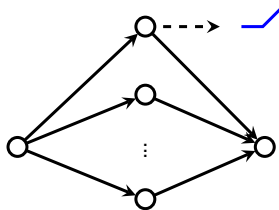
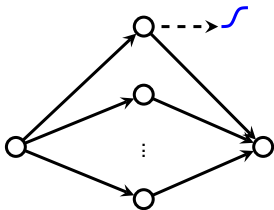
Squared Exponential kernel

$$(z; z^0) = \exp \left(-\frac{kz \cdot z^0 k^2}{2 \cdot 2} \right)$$

Matern- kernel

$$(z; z^0) = \frac{z^1}{(\quad)} \frac{p \sqrt{2}}{kz} z^0 k \quad K \quad \frac{p \sqrt{2}}{1} kz \quad z^0 k$$

Kernels



Reproducing Kernel Hilbert Space

RKHS:

$$H = \{ f(x) = \sum_{m=1}^P w_m \phi_m(x) \}$$

Inner product $\langle f, g \rangle = \sum_{k=1}^P w_f^k w_g^k$

$\|f\|_H = \|w\|$

$\phi_m = \frac{1}{\sqrt{m}}$ form an orthonormal basis

Kernel Based Regression

Provided a dataset of observation:

$$(z_j; Y(z_j))_{j=1}^t, Y(z_j) = f(z_j) + \epsilon_j$$

Regularized Least Squares Error:

$$\hat{f} = \arg \min_{g \in \mathcal{H}} \sum_{j=1}^t (Y(z_j) - g(z_j))^2 + \lambda \|g\|_{\mathcal{H}}^2$$

Kernel-Based Regression

Predictor:

$$\hat{f}_t(z) = \hat{\gamma}_t(z)(K_t + I)^{-1}y_t$$

$$\hat{\gamma}_t(z) = [k(z_1; z); k(z_2; z); \dots; k(z_t; z)]$$

$$K_t = [k(z_i; z_j)]_{i,j=1}^t$$

$$y_t = [Y(z_1); Y(z_2); \dots; Y(z_t)]$$

Kernel-Based Regression

Uncertainty estimator:

$$\hat{\sigma}_t^2(z) = \hat{\sigma}_t^2(z; z) \hat{\sigma}_t^2(z) (K_t + I)^{-1} \hat{\sigma}_t^2(z)$$

Kernel-Based Regression

Uncertainty estimator:

$$\sigma_t^2(z) = \mathbf{k}_t(z; z) - \mathbf{k}_t^T(z)(\mathbf{K}_t + \mathbf{I})^{-1} \mathbf{k}_t(z)$$

Closed form expressions for prediction and uncertainty quantification!

RL with kernel-based function approximation

IDEA: Approximate the Q-functions as a function in RKHS

RL with kernel-based function approximation

IDEA: Approximate the Q -functions as a function in RKHS

Want: Find an ϵ -optimal policy with a computational and sample complexity polynomial in $1/\epsilon$ and H

RL with kernel-based function approximation

IDEA: Approximate the Q -functions as a function in RKHS

Want: Find an ϵ -optimal policy with a computational and sample complexity polynomial in $1/\epsilon$ and H

[possibly some kernel parameters]

RL with kernel-based function approximation

IDEA: Approximate the Q -functions as a function in RKHS

Want: Find an ϵ -optimal policy with a computational and sample complexity polynomial in $1/\epsilon$ and H

|possibly some kernel parameters|

independently of \mathcal{S} and A .

RL with kernel-based function approximation

IDEA: Approximate the Q -functions as a function in RKHS

Want: Find an ϵ -optimal policy with a computational and sample complexity polynomial in $1/\epsilon$ and H

[possibly some kernel parameters]

independently of \mathcal{S} and A .

Sample complexity: $\mathcal{O}\left(\frac{1}{\epsilon}\right)^2$

RL with kernel-based function approximation

IDEA: Approximate the Q-functions as a function in RKHS

Want: Find an ϵ -optimal policy with a computational and sample complexity polynomial in $1/\epsilon$ and H

[possibly some kernel parameters]

independently of \mathcal{S} and A .

Sample complexity $\tilde{O}\left(\frac{1}{\epsilon}\right)$?

RL with kernel-based function approximation

RL with kernel-based function approximation

Effective Dimension:

$$\left| \begin{array}{c} 1; \dots; z_t; \dots; D; \\ \hline \text{D dimension} \end{array} \right|^{D+1}$$

$$D = \frac{1}{2} \log \det(I + \frac{1}{\lambda} K_t)$$

In the linear setting: $D = d$

For Squared Exponential kernel:

$$D = \text{polylog}(T)$$

For Matern kernel:

$$D = T^{\frac{d}{d+1}} \text{ [Vakili et al., 2021]}$$

RL with Kernel-based FA

Kernel-based transition assumption

For all s^0 : $P(s^0 | s; a) \geq H$

RL with Kernel-based FA

Kernel-based transition assumption

For all s^0 : $P(s^0 | s; a) \in \mathcal{H}$

A significant generalization of linear models

RL with Kernel-based FA

Kernel-based transition assumption

For all s^0 : $P(s^0|s; a) \geq H$

A significant generalization of linear models

Linear model is a special case with linear kernel:

$$(s; a; s^0, a^0) = \langle \phi(s; a), \phi(s^0, a^0) \rangle$$

RL with Kernel-based FA

Kernel-based transition assumption

For all s^0 : $P(s^0|s; a) \in H$

A significant generalization of linear models

Linear model is a special case with linear kernel:

$$(s; a; s^0, a^0) = \langle (s; a), (s^0, a^0) \rangle$$

RKHS of common kernels can approximate almost all continuous functions

RL with Kernel-based FA

Kernel-based transition assumption

For all s^0 : $P(s^0; s; a) \in H$

A significant generalization of linear models

Linear model is a special case with linear kernel:

$$(s; a; s^0, a^0) = \langle (s; a), (s^0, a^0) \rangle$$

RKHS of common kernels can approximate almost all continuous functions

For integrable $V : S \rightarrow \mathbb{R}$, $[PV] = \int_{s^0} P(s^0; s; a) V(s^0) \in H_k$

Optimistic approximate DP goes kernelized

IDEA: Combine the techniques for tabular MDPs with exploration bonuses borrowed from the linear kernel bandit literature!

UCB-VI [Azar et al., 2017]:

Backtrack $h = H; H-1; \dots; 1$: run optimistic value iteration

$$Q_h = r_h + \underbrace{P_h}_{\text{model estimate}} V_{h+1} + \underbrace{b_h}_{\text{exploration bonus}}$$

and set $V_h(s) = \max_a Q_h(s; a)$ for all $s; a$.

Forward $h = 1; 2; \dots; H$: take actions according to greedy policy

$$a_h(s) = \arg \max_a Q_h(s; a)$$

But how do we define P_h and b_h ?

Kernel-based optimistic value iteration (KOVI)

Transition model P_h can be defined **implicitly** via least-squares:

- | Solve the regularized linear regression problem

$$\hat{f}_h = \arg \min_{f \in \mathcal{H}} \sum_{t=1}^P (V_{h+1}(s_{h+1;t}) - f(s_{h;t}; a_{h;t}))^2 + \lambda \|f\|_{\mathcal{H}}^2$$

- | That provides a prediction

$$[P_h \backslash V_{h+1;k}](s; a) = \hat{f}_h(s; a) = \sum_{h;k} \phi_{h;k}(s; a) (K_{h;k} + I)^{-1} v_{h;k}$$

$$v_{h;k} = [V_{h+1}(s_{h+1;1}); V_{h+1}(s_{h+1;2}); \dots; V_{h+1}(s_{h+1;k})]$$

- | Also, an uncertainty quantification (variance)

$$\hat{\Sigma}_{h;k}(s; a) = \phi_{h;k}(s; a); \phi_{h;k}(s; a) \sum_{h;k} \phi_{h;k}(s; a) (K_{h;k} + I)^{-1} \phi_{h;k}(s; a)$$

Kernel-based optimistic value iteration (KOVI)

The prediction and variance give us an upper confidence bound on Q

$$Q_{h;k}(s; a) = r_h(s; a) + [P_h V_{h+1}](s; a) + \beta_{h;k}(s; a)$$

This is then used to compute an UCB Q as

$$V_{h;k}(s) = \max_a Q_{h;k}(s; a)$$

Performance guarantees

Theorem [Yang et al., 2020] The regret of KOVI satisfies
$$\text{Regret}(K) = O\left(H^{2p} \frac{D^2 + D \log N}{K}\right).$$

Performance guarantees

Theorem [Yang et al., 2020] The regret of KOVI satisfies

$$\text{Regret}(K) = \mathcal{O}\left(H^{2p} \sqrt{(D^2 + D \log N)K}\right).$$

Proof ideas:

Prove confidence bounds

$$|P_h \setminus V_{h+1;k}|(s; a) - |P_h V_{h+1;k}|(s; a)| \leq \epsilon_{h;k}(s; a):$$

Using standard techniques, show

$$\text{Regret}(K) \leq \sum_{h,k} \epsilon_{h;k}(s_{h;k}; a_{h;k}):$$

Kernelized elliptical potential lemma (e.g., Srinivas et al., 2010)

$$\sum_{h,k} \epsilon_{h;k}(s_{h;k}; a_{h;k}) \leq H^p \sqrt{K D \log(K)}:$$

Kernel-based Setting - Analysis

I We need a confidence bound of the form

$$\hat{f}(s; a) - [P_h V_{h+1; k}](s; a) \leq \epsilon_h(s; a):$$

Kernel-based Setting - Analysis

- I We need a confidence bound of the form

$$\hat{f}(s; a) - [P_h V_{h+1; k}](s; a) \leq \epsilon_n(s; a):$$

- I For a fixed $f \in H$ with non-adaptive inputs z_1, \dots, z_k ,

$$\epsilon_n \leq k f_{k_H} + \frac{H}{p} \sqrt{d \log\left(\frac{T}{\epsilon}\right)}$$

Kernel-based Setting - Analysis

I We need a confidence bound of the form

$$\hat{f}(s; a) - [P_h V_{h+1;k}](s; a) \leq \epsilon_n(s; a):$$

I For a fixed $f \in \mathcal{H}$ with non-adaptive inputs $z_1; \dots; z_k$,

$$\epsilon_n \leq k f_{k_H} + \frac{q}{p} \sqrt{d \log\left(\frac{T}{\epsilon}\right)}$$

Challenge 1: Inputs $(s_1; a_1); \dots; (s_k; a_k)$ are adaptive!

Kernel-based Setting - Analysis

I We need a confidence bound of the form

$$\hat{f}(s; a) - [P_h V_{h+1; k}](s; a) \leq \left(\frac{1}{n} \right) h(s; a):$$

I For a fixed $f \in \mathcal{H}$ with non-adaptive inputs $z_1; \dots; z_k$,

$$\left(\frac{1}{n} \right) k f k_H + \frac{1}{\epsilon} \sqrt{q \log\left(\frac{1}{\epsilon}\right)}$$

Challenge 1: Inputs $(s_1; a_1); \dots; (s_k; a_k)$ are adaptive!

Solution: Self-normalized concentration inequalities for vector-valued martingales extended to kernel setting [[Abbasi-Yadkori, 2013](#), [Whitehouse et al., 2023](#)]:

$$\left(\frac{1}{n} \right) k f k_H + \frac{1}{\epsilon} \sqrt{q \log\left(\frac{1}{\epsilon}\right)}$$

Kernel-based Setting - Analysis

I We need a confidence bound of the form

$$\hat{f}(s; a) - [P_h V_{h+1; k}](s; a) \leq \epsilon_n(s; a):$$

I For a fixed $f \in \mathcal{H}$ with non-adaptive inputs z_1, \dots, z_k ,

$$\epsilon_n \leq k f_{k_H} + \frac{q}{\epsilon} \sqrt{d \log\left(\frac{1}{\epsilon}\right)}$$

Challenge 1: Inputs $(s_1; a_1); \dots; (s_k; a_k)$ are adaptive!

Solution: Self-normalized concentration inequalities for vector-valued martingales extended to kernel setting [[Abbasi-Yadkori, 2013](#), [Whitehouse et al., 2023](#)]:

$$\epsilon_n \leq k f_{k_H} + \frac{q}{\epsilon} \sqrt{D + \log\left(\frac{1}{\epsilon}\right)}$$

Challenge 2: $f = P_h V_{h+1; k}$ is **not** fixed, but depends on past data!

Kernel-based Setting - Analysis

I We need a confidence bound of the form

$$\hat{f}(s; a) - [P_h V_{h+1; k}](s; a) \leq \epsilon_n(s; a):$$

I For a fixed $f \in \mathcal{H}$ with non-adaptive inputs z_1, \dots, z_k ,

$$\epsilon_n \leq \epsilon_{k_H} + \frac{q}{p} \sqrt{d \log\left(\frac{1}{\delta}\right)}$$

Challenge 1: Inputs $(s_1; a_1); \dots; (s_k; a_k)$ are adaptive!

Solution: Self-normalized concentration inequalities for vector-valued martingales extended to kernel setting [[Abbasi-Yadkori, 2013](#), [Whitehouse et al., 2023](#)]:

$$\epsilon_n \leq \epsilon_{k_H} + \frac{q}{p} \sqrt{D + \log\left(\frac{1}{\delta}\right)}$$

Challenge 2: $f = P_h V_{h+1; k}$ is **not** fixed, but depends on past data!

Solution: Covering number argument

Covering Number Argument

- I Notice that all value functions $V_{h,k}$ belong to the function class

$$V = \{ V(s) = \min_{f \in H} \max_a f(s; a) + \epsilon \}$$

- I How many functions V are required to cover up to ϵ error?

Covering Number Argument

- | Notice that all value functions $V_{h,k}$ belong to the function class

$$V = \{ V(s) = \min_{f \in H} \max_a f(s; a) + \epsilon \}$$

- | How many functions V are required to cover V up to ϵ error?

Covering Number Argument

- | Notice that all value functions $V_{h,k}$ belong to the function class

$$V = \{ V(s) = \min_{f \in H} \max_a f(s; a) + \epsilon \}$$

- | How many functions V are required to cover V up to ϵ error?

Covering Number Argument

I We can now use a union-bound argument

$$\left(\right) = \left(=N \right) k f k_{H_k} + p^H \frac{q}{D + \log N + 1}$$

Covering Number Argument

- I We can now use a union-bound argument

$$\left(\frac{1}{K} \right) = \frac{1}{K} \left(\sum_{k=1}^K \frac{1}{K} \right) \leq \frac{1}{K} \left(\sum_{k=1}^K \frac{1}{K} \right) + \frac{1}{K} \frac{H^p}{D^2 K + D \log N K}$$

- I Regret $\left(\frac{1}{K} \right)$ [Yang et al., 2020]

$$\text{Regret}(K) = O\left(H^2 \frac{D^2 K + D \log N K}{K}\right)$$

Covering Number Argument

- | We can now use a union-bound argument

$$f_{H_k} = \frac{H}{D + \log N} + \frac{1}{K}$$

- | Regret $\left(\frac{1}{K}\right)$ [Yang et al., 2020]

$$\text{Regret}(K) = \tilde{O}\left(H^2 \sqrt{D^2 K + D \log N K}\right)$$

- | Sample Complexity

- Very smooth kernels D and $\log(N)$ $\text{poly} \log(K)$

$$\tilde{O}\left(\frac{1}{2}\right)$$

- In general could be vacuous!

Optimistic Closure

Chowdhury and Oliveira [2023] Optimistic Closure Assumption:

$V \leq H$

$$V = V(s) = \min\{H, \max_a \hat{f}(s, a) + V(s, a)\}$$

Optimistic Closure

Chowdhury and Oliveira [2023] Optimistic Closure Assumption:

$$V \leq H$$

$$V = V(s) = \min\{H, \max_a \hat{f}(s, a) + \gamma V(s, a)\}$$

Idea: Leverage kernel mean embedding

$$\text{Regret}(K) = \tilde{O}(H^2 D \sqrt{K})$$

Optimistic Closure

Chowdhury and Oliveira [2023] Optimistic Closure Assumption:

$$V \leq H$$

$$V = \min_{s,a} \{H, \max_a \hat{f}(s, a) + \dots (s, a)\}$$

Idea: Leverage kernel mean embedding

$$\text{Regret}(K) = \tilde{O}(H^2 D \overline{K})$$

Doen not hold in the linear setting!

Open Problem Vakili [2024]

- (a) Can a **no-regret** learning algorithm be designed?
- (b) What is the minimum regret growth rate with K (and also H)? And, can a learning algorithm be designed to achieve order optimal (or near-optimal) regret performance, closely aligning with the established lower bound?

Some References

- Chowdhury and Gopalan [2019]
- Yang et al. [2020]
- Domingues et al. [2021]
- Vakili and Olkhovskaya [2023]
- ...

References I

- Y. Abbasi-Yadkori. Online learning for linearly parametrized control problems. *PhD Thesis, University of Alberta*, 2013.
- Y. Abbasi-Yadkori, D. Pál, and C. Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in Neural Information Processing Systems*, 24, 2011.
- A. Ayoub, Z. Jia, C. Szepesvari, M. Wang, and L. Yang. Model-based reinforcement learning with value-targeted regression. In *International Conference on Machine Learning*, pages 463–474. PMLR, 2020.
- M. G. Azar, I. Osband, and R. Munos. Minimax regret bounds for reinforcement learning. In *International conference on machine learning*, pages 263–272. PMLR, 2017.
- S. R. Chowdhury and A. Gopalan. Online learning in kernelized Markov decision processes. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 3197–3205. PMLR, 2019.
- S. R. Chowdhury and R. Oliveira. Value function approximations via kernel embeddings for no-regret reinforcement learning. In *Asian Conference on Machine Learning*, pages 249–264. PMLR, 2023.
- O. D. Domingues, P. Ménard, M. Pirotta, E. Kaufmann, and M. Valko. Kernel-based reinforcement learning: A finite-time analysis. In *International Conference on Machine Learning*, pages 2783–2792. PMLR, 2021.

References II

- S. Du, S. Kakade, J. Lee, S. Lovett, G. Mahajan, W. Sun, and R. Wang. Bilinear classes: A structural framework for provable generalization in RL. In *International Conference on Machine Learning*, pages 2826–2836. PMLR, 2021.
- C. Jin, Z. Allen-Zhu, S. Bubeck, and M. I. Jordan. Is Q-learning provably efficient? *Advances in neural information processing systems*, 31, 2018.
- C. Jin, Z. Yang, Z. Wang, and M. I. Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143. PMLR, 2020.
- C. Jin, Q. Liu, and S. Miryoosefi. Bellman eluder dimension: New rich classes of RL problems, and sample-efficient algorithms. *Advances in neural information processing systems*, 34:13406–13418, 2021.
- C. Jin, Z. Yang, Z. Wang, and M. I. Jordan. Provably efficient reinforcement learning with linear function approximation. *Mathematics of Operations Research*, 48(3):1496–1521, 2023.
- T. Lattimore and C. Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- J. Lee and M.-h. Oh. Demystifying linear mdps and novel dynamics aggregation framework. In *The Twelfth International Conference on Learning Representations*, 2024.

References III

- A. Modi, J. Chen, A. Krishnamurthy, N. Jiang, and A. Agarwal. Model-free representation learning and exploration in low-rank mdps. *Journal of Machine Learning Research*, 25(6):1–76, 2024.
- A. Moulin and G. Neu. Optimistic planning by regularized dynamic programming. In *International Conference on Machine Learning*, pages 25337–25357. PMLR, 2023.
- G. Neu and C. Pike-Burke. A unifying view of optimism in episodic reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1392–1403, 2020.
- N. Srinivas, A. Krause, S. Kakade, and M. Seeger. Gaussian process optimization in the bandit setting: no regret and experimental design. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, pages 1015–1022, 2010.
- S. Vakili. Open problem: Order optimal regret bounds for kernel-based reinforcement learning. In *The Thirty Seventh Annual Conference on Learning Theory*, pages 5340–5344. PMLR, 2024.
- S. Vakili and J. Olkhovskaya. Kernelized reinforcement learning with order optimal regret bounds. *Advances in Neural Information Processing Systems*, 36, 2023.
- S. Vakili, K. Khezeli, and V. Picheny. On information gain and regret bounds in Gaussian process bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 82–90. PMLR, 2021.

References IV

- G. Weisz, P. Amortila, and C. Szepesvári. Exponential lower bounds for planning in mdps with linearly-realizable optimal action-value functions. In *Algorithmic Learning Theory*, pages 1237–1264. PMLR, 2021.
- G. Weisz, A. György, T. Kozuno, and C. Szepesvári. Confident approximate policy iteration for efficient local planning in q -realizable mdps. *Advances in Neural Information Processing Systems*, 35:25547–25559, 2022.
- G. Weisz, A. György, and C. Szepesvári. Online RL in linearly q -realizable MDPs is as easy as in linear MDPs if you learn what to ignore. *Advances in Neural Information Processing Systems*, 36, 2023.
- J. Whitehouse, A. Ramdas, and S. Z. Wu. On the sublinear regret of gp-ucb. *Advances in Neural Information Processing Systems*, 36, 2023.
- L. Yang and M. Wang. Sample-optimal parametric q-learning using linearly additive features. In *International conference on machine learning*, pages 6995–7004. PMLR, 2019.
- L. Yang and M. Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. In *International Conference on Machine Learning*, pages 10746–10756. PMLR, 2020.

References V

- Z. Yang, C. Jin, Z. Wang, M. Wang, and M. Jordan. Provably efficient reinforcement learning with kernel and neural function approximations. *Advances in Neural Information Processing Systems*, 33:13903–13916, 2020.
- A. Zanette, D. Brandfonbrener, E. Brunskill, M. Pirotta, and A. Lazaric. Frequentist regret bounds for randomized least-squares value iteration. In *International Conference on Artificial Intelligence and Statistics*, pages 1954–1964. PMLR, 2020a.
- A. Zanette, A. Lazaric, M. Kochenderfer, and E. Brunskill. Learning near optimal policies with low inherent Bellman error. In H. D. III and A. Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 10978–10989. PMLR, 13–18 Jul 2020b.
- D. Zhou, J. He, and Q. Gu. Provably efficient reinforcement learning for discounted mdps with feature mapping. In *International Conference on Machine Learning*, pages 12793–12802. PMLR, 2021.