

# FoSML: Lecture 21

The previous result about approximate kernels requires understanding covering numbers. In particular,  $N(2R, r)$  appears as a constant driving our bound. In general,  $N(2R, r)$  depends on the geometry and dimension of  $X$ .

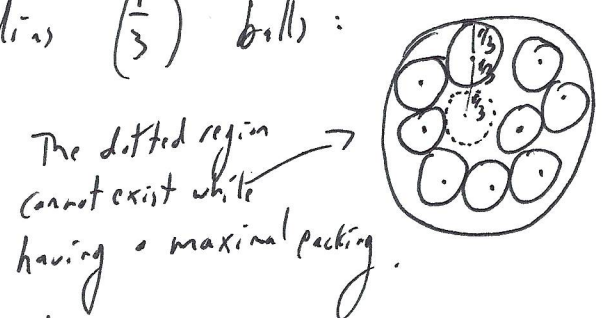
In the case  $X \subseteq \mathbb{R}^D$ , it can be estimated as follows:

Lemma: Let  $X \subseteq \mathbb{R}^D$  be compact and let  $R = \min_{\delta > 0} \{ \delta \mid X \subset B_\delta(0) \}$ .

Then  $N(R, r) \leq \left(\frac{3R}{r}\right)^D$ .

Proof: Up to constants independent of  $\delta$ ,  $\text{Vol}_D(B_\delta(0)) = \delta^D$ . Then there is no way to fit more than  ~~$R^D$~~   $R^D / (r/3)^D = (3R/r)^D$  balls of radius  $r/3$  in  $B_R(0)$  without the balls intersecting; indeed, achieving this bound would require no "wasted space" in the packing.

Now, consider a maximal packing of at most  $(3R/r)^D$  balls of radius  $r/3$  into the ball  $B_R(0)$ . Then every element of  $B_R(0)$  is distance at most  $r$  from the center of one of the radius  $(\frac{r}{3})$  balls:



Then if the balls of radius  $\frac{r}{3}$  are inflated to radius  $r$ , they necessarily cover  $B_R(0)$ . Thus, we have a cover of radius  $r$  balls of cardinality  $\left(\frac{3R}{r}\right)^D \Rightarrow N(R, r) \leq \left(\frac{3R}{r}\right)^D$ . ■

Combining these estimates yields:

②

Theorem: Let  $K$  be a continuously differentiable kernel function that satisfies our conditions (shift-invariant, associated non-negative measure  $\mu$ ). Suppose  $\sigma_\mu^2 = \mathbb{E}_{w \sim \mu} (\|w\|_2^2) < \infty$  and  $X \subset \mathbb{R}^D$  is compact and contained in  $B_R(0)$ .

Then for  $\Psi: X \rightarrow \mathbb{R}^{2D}$  defined as before, ~~then~~ if  $\epsilon \in (0, 32R\sigma_\mu)$ ,

$$P \left( \sup_{x, x' \in X} |\langle \Psi(x), \Psi(x') \rangle - K(x, x')| \geq \epsilon \right)$$

$$\leq \left( \frac{48R\sigma_\mu}{\epsilon} \right)^2 \exp \left( \frac{-D\epsilon^2}{4(D+2)} \right)$$

$$\text{Proof: The result is immediate from setting } r = \left( \frac{2(6R)^D \exp \left( -\frac{D\epsilon^2}{8} \right)}{\left( \frac{4\sigma_\mu}{\epsilon} \right)^2} \right)^{\frac{2}{D+2}}$$

in the ~~then~~ lemma.  $\blacksquare$

Reinforcement learning: planning & learning when an agent wants to maximize some well-defined objective.

Instead of getting a random sample of training data  $\{(x_i, y_i)\}_{i=1}^n$ , an agent collects information about its environment through a sequence of actions.

After engaging in an action, an agent receives two pieces of information:

- (a) Its current state in the environment
- (b) A  $\mathbb{R}$ -valued reward, which is goal-dependent.

The agent aims to formulate a course of action, or policy, to maximize its

rewards. This policy setting problem needs to manage the exploration-exploitation (3) trade-off, i.e. should the agent exploit the rewards already known, or should the agent risk what is known to explore new possibilities for increased rewards.

. If the environment is known, we have a planning problem. If the environment is unknown, we have a learning problem.

- A classic framework for RL is Markov decision processes (MDP):

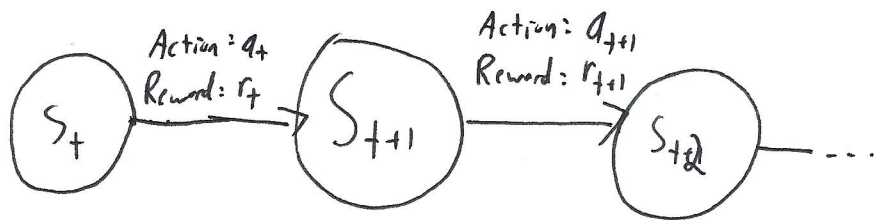
Defn: A MDP is defined by:

- (a.) A set of states  $S$ .
- (b.) An initial state  $s_0 \in S$ .
- (c.) A set of actions  $A$ .
- (d.) A transition probability <sup>dist.</sup> over destination states  $\delta(s,a) = P(\overset{s'}{s'} = \delta(s,a) \mid s,a)$
- (e.) A reward probability distribution over reward states  $r(s,a)$ :  
 $P(r' = r(s,a) \mid s,a)$ .

. Since the destination and reward probabilities depend only on the current state and action  $(s,a)$ , and not previous states or actions, the process is Markovian.

. We will think of taking discrete time steps, though MDP can be adapted to allow for continuous time parametrization. Let  $\{0, \dots, T\}$  denote the decision epochs with finite time horizon  $T$ .

. If  $|S|, |A| < \infty$ , we say the MDP is finite.



Action: agent chooses  
 Reward, next state  
 stochastic, depending on  
 action and current state

(4)

Defn: For a space of actions  $A$ , let  $\mathcal{P}(A) = \{\text{probability measures on } A\}$ .

Defn: Given an MDP with state space  $S$  and action space  $A$ , a policy is a function  $\pi: S \rightarrow \mathcal{P}(A)$ . A policy is deterministic if  $|\text{supp}(\pi(s))| = 1, \forall s \in S$ .

Remark: This definition of policy is time-independent; one could consider a time-dependent variant with  $\pi_t: S \rightarrow \mathcal{P}(A)$  changing with  $t$ . We say a policy that doesn't depend on time is stationary.

• What's a good policy? One that maximizes rewards across time:

ex: If  $T < \infty$  and each time is equally valuable, we may want  $\pi$  to maximize

$$\sum_{t=0}^T r(s_t, \pi(s_t))$$

state ~~the~~ distribution over actions, given that the agent is at state  $s_t$  at time  $t$ .

ex: If  $T = +\infty$ , it is necessary to discount in time. This can be done geometrically, for example, by considering  $\gamma \in (0, 1)$  and choosing  $\pi$  to maximize

$$\sum_{t=0}^{\infty} \gamma^t \cdot r(s_t, \pi(s_t))$$

discounting in time.

Defn: A policy  $\pi$  for an MDP at state  $s$  <sup>has value</sup>  $V_\pi(s)$  defined as (5)

(a.)  $T < \infty$ :  $\mathbb{E}_{a_t \sim \pi(s_t)} \left( \sum_{t=0}^T r(s_t, a_t) \mid s_0 = s \right)$

(b.)  $T = +\infty$ , discounting factor  $\gamma$ :  $\mathbb{E}_{a_t \sim \pi(s_t)} \left( \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right)$

So, we want to make  $V_\pi(s)$  large, for a given  $s$ , or perhaps averaged over all  $s \in \mathcal{S}$ .

Defn: A policy  $\pi^*$  for an MDP is optimal if  $\forall s \in \mathcal{S}$  and  $\forall \pi = s \rightarrow \mathcal{P}(A)$ ,  $V_{\pi^*}(s) \geq V_\pi(s)$ , i.e.  $\pi^*$  yields maximal reward for all initial states.

$\square$ : When does such a  $\pi^*$  exist? It appears optimality is a hard condition. We will show such states exist, as long as  $|\mathcal{A}|, |\mathcal{S}| < \infty$ .

Defn: ~~...~~ The state-action function  $Q$  for a policy  $\pi$  is

$$Q_\pi(s, a) = \mathbb{E}(r(s, a)) + \mathbb{E}_{a_t \sim \pi(s_t)} \left( \sum_{t=1}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right)$$
$$= \mathbb{E}(r(s, a) + \gamma V_\pi(s_1) \mid s_0 = s, a_0 = a).$$

Note that  $\mathbb{E}_{a \sim \pi(s)} (Q_\pi(s, a)) = V_\pi(s)$ .

Theorem (policy improvement): Let  $\pi, \pi'$  be policies for a common MDP. (6)

$$\text{Then: } \left\{ \forall s \in \mathcal{S}, \mathbb{E}_{a \sim \pi'(s)} (Q_{\pi}(s, a)) \geq \mathbb{E}_{a \sim \pi(s)} (Q_{\pi}(s, a)) \right\} \quad (a.)$$

$\Downarrow$

$$\left\{ \forall s \in \mathcal{S}, V_{\pi'}(s) \geq V_{\pi}(s) \right\}. \quad (b.)$$

Moreover, strict inequality for at least one  $s \in \mathcal{S}$  on LHS  $\Rightarrow$  strict inequality for at least one  $s \in \mathcal{S}$  on RHS.

Proof: Suppose  $\pi, \pi'$  satisfy (a). Then,  $\forall s \in \mathcal{S}$ ,

$$V_{\pi}(s) = \mathbb{E}_{a \sim \pi(s)} (Q_{\pi}(s, a))$$

$$\leq \mathbb{E}_{a \sim \pi'(s)} (Q_{\pi}(s, a))$$

$$= \mathbb{E}_{a \sim \pi'(s)} (r(s, a) + \gamma V_{\pi}(s_1) \mid s_0 = s)$$

$$= \mathbb{E}_{a \sim \pi'(s)} (r(s, a) + \gamma \mathbb{E}_{a_1 \sim \pi(s_1)} (Q_{\pi}(s_1, a_1)) \mid s_0 = s)$$

$$\leq \mathbb{E}_{a \sim \pi'(s)} (r(s, a) + \gamma \mathbb{E}_{a_1 \sim \pi'(s_1)} (Q_{\pi}(s_1, a_1)) \mid s_0 = s)$$

$$= \mathbb{E}_{\substack{a \sim \pi'(s) \\ a_1 \sim \pi'(s_1)}} (r(s, a) + \gamma r(s_1, a_1) + \gamma^2 V_{\pi}(s_2) \mid s_0 = s).$$

Proceeding inductively, for any  $T \geq 1$ ,

$$V_{\pi}(s) \leq \mathbb{E}_{a_t \sim \pi'(s_t)} \left( \sum_{t=0}^T \gamma^t \mathbb{E} (r(s_t, a_t)) + \gamma^{T+1} V_{\pi}(s_{T+1}) \mid s_0 = s \right)$$

Taking  $T \rightarrow \infty$  gives a convergent sum equal to  $V_{\pi}(s)$ . ■