

---

## Appendix: Lipschitz Continuity in Model-based Reinforcement Learning

We first restate the core lemmas, theorems, and claims presented in our paper below:

**Lemma 1.** A generalized transition function  $\widehat{T}_{\mathcal{G}}$  induced by a Lipschitz model class  $F_g$  is Lipschitz with a constant:

$$K_{W,W}^A(\widehat{T}_{\mathcal{G}}) := \sup_a \sup_{\mu_1, \mu_2} \frac{W(\widehat{T}_{\mathcal{G}}(\cdot | \mu_1, a), \widehat{T}_{\mathcal{G}}(\cdot | \mu_2, a))}{W(\mu_1, \mu_2)} \leq K_F$$

**Lemma 2.** (Composition Lemma) Define three metric spaces  $(M_1, d_1)$ ,  $(M_2, d_2)$ , and  $(M_3, d_3)$ . Define Lipschitz functions  $f : M_2 \mapsto M_3$  and  $g : M_1 \mapsto M_2$  with constants  $K_{d_2, d_3}(f)$  and  $K_{d_1, d_2}(g)$ . Then,  $h : f \circ g : M_1 \mapsto M_3$  is Lipschitz with constant  $K_{d_1, d_3}(h) \leq K_{d_2, d_3}(f)K_{d_1, d_2}(g)$ .

**Theorem 1.** Define a  $\Delta$ -accurate  $\widehat{T}_{\mathcal{G}}$  with the Lipschitz constant  $K_F$  and an MDP with a Lipschitz transition function  $T_{\mathcal{G}}$  with constant  $K_T$ . Let  $\bar{K} = \min\{K_F, K_T\}$ . Then  $\forall n \geq 1$ :

$$\delta(n) := W(\widehat{T}_{\mathcal{G}}^n(\cdot | \mu), T_{\mathcal{G}}^n(\cdot | \mu)) \leq \Delta \sum_{i=0}^{n-1} (\bar{K})^i.$$

**Theorem 2.** Assume a Lipschitz model class  $F_g$  with a  $\Delta$ -accurate  $\widehat{T}$  with  $\bar{K} = \min\{K_F, K_T\}$ . Further, assume a Lipschitz reward function with constant  $K_R = K_{d_{\mathcal{S}}, \mathbb{R}}(R)$ . Then  $\forall s \in \mathcal{S}$  and  $\bar{K} \in [0, \frac{1}{\gamma})$

$$|V_T(s) - V_{\widehat{T}}(s)| \leq \frac{\gamma K_R \Delta}{(1 - \gamma)(1 - \gamma \bar{K})}.$$

**Lemma 3.** Given a Lipschitz function  $f : \mathcal{S} \mapsto \mathbb{R}$  with constant  $K_{d_{\mathcal{S}}, d_{\mathbb{R}}}(f)$ :

$$K_{d_{\mathcal{S}}, d_{\mathbb{R}}}^A \left( \int \widehat{T}(s' | s, a) f(s') ds' \right) \leq K_{d_{\mathcal{S}}, d_{\mathbb{R}}}(f) K_{d_{\mathcal{S}}, W}^A(\widehat{T}).$$

**Lemma 4.** The following operators (Asadi & Littman, 2017) are Lipschitz with constants:

1.  $K_{\|\cdot\|_{\infty}, d_R}(\max(x)) = K_{\|\cdot\|_{\infty}, d_R}(\text{mean}(x)) = K_{\|\cdot\|_{\infty}, d_R}(\epsilon\text{-greedy}(x)) = 1$
2.  $K_{\|\cdot\|_{\infty}, d_R}(mm_{\beta}(x)) := \frac{\log \frac{\sum_i e^{\beta x_i}}{n}}{\beta} = 1$
3.  $K_{\|\cdot\|_{\infty}, d_R}(\text{boltz}_{\beta}(x)) := \frac{\sum_{i=1}^n x_i e^{\beta x_i}}{\sum_{i=1}^n e^{\beta x_i}} \leq \sqrt{|A|} + \beta V_{\max}|A|$

**Theorem 3.** For any non-expansion backup operator  $f$  outlined in Lemma 4, GVI computes a value function with a Lipschitz constant bounded by  $\frac{K_{d_{\mathcal{S}}, d_{\mathbb{R}}}^A(R)}{1 - \gamma K_{d_{\mathcal{S}}, W}^A(T)}$  if  $\gamma K_{d_{\mathcal{S}}, W}^A(T) < 1$ .

We now provide proofs of various results mentioned in the paper:

**Claim 1.** In a finite MDP, transition probabilities can be expressed using a finite set of deterministic functions and a distribution over the functions.

*Proof.* Let  $Pr(s, a, s')$  denote the probability of a transition from  $s$  to  $s'$  when executing the action  $a$ . Define an ordering over states  $s_1, \dots, s_n$  with an additional unreachable state  $s_0$ . Now define the cumulative probability distribution:

$$C(s, a, s_i) := \sum_{j=0}^i Pr(s, a, s_j).$$

Further define  $L$  as the set of distinct entries in  $C$ :

$$L := \left\{ C(s, a, s_i) \mid s \in \mathcal{S}, i \in [0, n] \right\}.$$

Note that, since the MDP is assumed to be finite, then  $|L|$  is finite. We sort the values of  $L$  and denote, by  $c_i$ ,  $i$ th smallest value of the set. Note that  $c_0 = 0$  and  $c_{|L|} = 1$ . We now build deterministic set of functions  $f_1, \dots, f_{|L|}$  as follows:  $\forall i = 1$  to  $|L|$  and  $\forall j = 1$  to  $n$ , define  $f_i(s) = s_j$  if and only if:

$$C(s, a, s_{j-1}) < c_i \leq C(s, a, s_j) .$$

We also define the probability distribution  $g$  over  $f$  as follows:

$$g(f_i|a) := c_i - c_{i-1} .$$

Given the functions  $f_1, \dots, f_{|L|}$  and the distribution  $g$ , we can now compute the probability of a transition to  $s_j$  from  $s$  after executing action  $a$ :

$$\begin{aligned} & \sum_i \mathbb{1}(f_i(s) = s_j) g(f_i|a) \\ &= \sum_i \mathbb{1}(C(s, a, s_{j-1}) < c_i \leq C(s, a, s_j)) (c_i - c_{i-1}) \\ &= C(s, a, s_j) - C(s, a, s_{j-1}) \\ &= Pr(s, a, s_j) , \end{aligned}$$

where  $\mathbb{1}$  is a binary function that outputs one if and only if its condition holds. We reconstructed the transition probabilities using distribution  $g$  and deterministic functions  $f_1, \dots, f_{|L|}$ .  $\square$

**Claim 2.** *Given a deterministic and linear transition model, and a linear reward signal, the bounds provided in Theorems 1 and 2 are both tight.*

Assume a linear transition function  $T$  defined as:

$$T(s) = Ks$$

Assume our learned transition function  $\hat{T}$ :

$$\hat{T}(s) := Ks + \Delta$$

Note that:

$$\max_s |T(s) - \hat{T}(s)| = \Delta$$

and that:

$$\min\{K_T, K_{\hat{T}}\} = K$$

First observe that the bound in Theorem 2 is tight for  $n = 2$ :

$$\forall s \quad \left| T(T(s)) - \hat{T}(\hat{T}(s)) \right| = \left| K^2s - K^2s + \Delta(1 + K) \right| = \Delta \sum_{i=0}^1 K^i$$

and more generally and after  $n$  compositions of the models, denoted by  $T^n$  and  $\hat{T}^n$ , the following equality holds:

$$\forall s \quad \left| T^n(s) - \hat{T}^n(s) \right| = \Delta \sum_{i=0}^{n-1} K^i$$

Lets further assume that the reward is linear:

$$R(s) = K_R s$$

Consider the state  $s = 0$ . Note that clearly  $v(0) = 0$ . We now compute the value predicted using  $\hat{T}$ , denoted by  $\hat{v}(0)$ :

$$\begin{aligned} \hat{v}(0) &= R(0) + \gamma R(0 + \Delta \sum_{i=0}^0 K^i) + \gamma^2 R(0 + \Delta \sum_{i=0}^1 K^i) + \gamma^3 R(0 + \Delta \sum_{i=0}^2 K^i) + \dots \\ &= 0 + \gamma K_R \Delta \sum_{i=0}^0 K^i + \gamma^2 K_R \Delta \sum_{i=0}^1 K^i + \gamma^3 K_R \Delta \sum_{i=0}^2 K^i + \dots \end{aligned}$$

$$= \gamma K_R \Delta \sum_{n=0}^{\infty} \gamma^n \sum_{i=0}^{n-1} K^i = \frac{\gamma K_R \Delta}{(1-\gamma)(1-\gamma \bar{K})},$$

and so:

$$|v(0) - \hat{v}(0)| = \frac{\gamma K_R \Delta}{(1-\gamma)(1-\gamma \bar{K})}$$

Note that this exactly matches the bound derived in our Theorem 2.

**Lemma 1.** A generalized transition function  $\widehat{T}_{\mathcal{G}}$  induced by a Lipschitz model class  $F_{\mathcal{G}}$  is Lipschitz with a constant:

$$K_{W,W}^A(\widehat{T}_{\mathcal{G}}) := \sup_a \sup_{\mu_1, \mu_2} \frac{W(\widehat{T}_{\mathcal{G}}(\cdot | \mu_1, a), \widehat{T}_{\mathcal{G}}(\cdot | \mu_2, a))}{W(\mu_1, \mu_2)} \leq K_F$$

*Proof.*

$$\begin{aligned} W(\widehat{T}(\cdot | \mu_1, a), \widehat{T}(\cdot | \mu_2, a)) &:= \inf_j \int_{s'_1} \int_{s'_2} j(s'_1, s'_2) d(s'_1, s'_2) ds'_1 ds'_2 \\ &= \inf_j \int_{s_1} \int_{s_2} \int_{s'_1} \int_{s'_2} \sum_f \mathbb{1}(f(s_1) = s'_1 \wedge f(s_2) = s'_2) j(s_1, s_2, f) d(s'_1, s'_2) ds'_1 ds'_2 ds_1 ds_2 \\ &= \inf_j \int_{s_1} \int_{s_2} \sum_f j(s_1, s_2, f) d(f(s_1), f(s_2)) ds_1 ds_2 \\ &\leq K_F \inf_j \int_{s_1} \int_{s_2} \sum_f g(f|a) j(s_1, s_2) d(s_1, s_2) ds_1 ds_2 \\ &= K_F \sum_f g(f|a) \inf_j \int_{s_1} \int_{s_2} j(s_1, s_2) d(s_1, s_2) ds_1 ds_2 \\ &= K_F \sum_f g(f|a) W(\mu_1, \mu_2) = K_F W(\mu_1, \mu_2) \end{aligned}$$

Dividing by  $W(\mu_1, \mu_2)$  and taking sup over  $a, \mu_1$ , and  $\mu_2$ , we conclude:

$$K_{W,W}^A(\widehat{T}) = \sup_a \sup_{\mu_1, \mu_2} \frac{W(\widehat{T}(\cdot | \mu_1, a), \widehat{T}(\cdot | \mu_2, a))}{W(\mu_1, \mu_2)} \leq K_F.$$

We can also prove this using the Kantorovich-Rubinstein duality theorem:

For every  $\mu_1, \mu_2$ , and  $a \in \mathcal{A}$  we have:

$$\begin{aligned} W(\widehat{T}_{\mathcal{G}}(\cdot | \mu_1, a), \widehat{T}_{\mathcal{G}}(\cdot | \mu_2, a)) &= \sup_{f: K_{d_{\mathcal{S}, \mathbb{R}}}(f) \leq 1} \int_{\mathcal{S}} (\widehat{T}_{\mathcal{G}}(s | \mu_1, a) - \widehat{T}_{\mathcal{G}}(s | \mu_2, a)) f(s) ds \\ &= \sup_{f: K_{d_{\mathcal{S}, \mathbb{R}}}(f) \leq 1} \int_{\mathcal{S}} \int_{s_0} (\widehat{T}(s | s_0, a) \mu_1(s_0) - \widehat{T}(s | s_0, a) \mu_2(s_0)) f(s) ds ds_0 \\ &= \sup_{f: K_{d_{\mathcal{S}, \mathbb{R}}}(f) \leq 1} \int_{\mathcal{S}} \int_{s_0} \widehat{T}(s | s_0, a) (\mu_1(s_0) - \mu_2(s_0)) f(s) ds ds_0 \\ &= \sup_{f: K_{d_{\mathcal{S}, \mathbb{R}}}(f) \leq 1} \int_{\mathcal{S}} \int_{s_0} \sum_t g(t | a) \mathbb{1}(t(s_0) = s) (\mu_1(s_0) - \mu_2(s_0)) f(s) ds ds_0 \\ &= \sup_{f: K_{d_{\mathcal{S}, \mathbb{R}}}(f) \leq 1} \sum_t g(t | a) \int_{s_0} \int_{\mathcal{S}} \mathbb{1}(t(s_0) = s) (\mu_1(s_0) - \mu_2(s_0)) f(s) ds ds_0 \\ &= \sup_{f: K_{d_{\mathcal{S}, \mathbb{R}}}(f) \leq 1} \sum_t g(t | a) \int_{s_0} (\mu_1(s_0) - \mu_2(s_0)) f(t(s_0)) ds_0 \end{aligned}$$

$$\begin{aligned}
&\leq \sum_t g(t | a) \sup_{f:K_{d_S, \mathbb{R}}(f) \leq 1} \int_{s_0} (\mu_1(s_0) - \mu_2(s_0)) f(t(s_0)) ds_0 \\
&\quad \text{composition of } f, t \text{ is Lipschitz with constant upper bounded by } K_F. \\
&= K_F \sum_t g(t | a) \sup_{f:K_{d_S, \mathbb{R}}(f) \leq 1} \int_{s_0} (\mu_1(s_0) - \mu_2(s_0)) \frac{f(t(s_0))}{K_F} ds_0 \\
&\leq K_F \sum_t g(t | a) \sup_{h:K_{d_S, \mathbb{R}}(h) \leq 1} \int_{s_0} (\mu_1(s_0) - \mu_2(s_0)) h(s_0) ds_0 \\
&= K_F \sum_t g(t | a) W(\mu_1, \mu_2) = K_F W(\mu_1, \mu_2)
\end{aligned}$$

Again we conclude by dividing by  $W(\mu_1, \mu_2)$  and taking sup over  $a, \mu_1$ , and  $\mu_2$ .  $\square$

**Lemma 2.** (Composition Lemma) Define three metric spaces  $(M_1, d_1)$ ,  $(M_2, d_2)$ , and  $(M_3, d_3)$ . Define Lipschitz functions  $f : M_2 \mapsto M_3$  and  $g : M_1 \mapsto M_2$  with constants  $K_{d_2, d_3}(f)$  and  $K_{d_1, d_2}(g)$ . Then,  $h : f \circ g : M_1 \mapsto M_3$  is Lipschitz with constant  $K_{d_1, d_3}(h) \leq K_{d_2, d_3}(f)K_{d_1, d_2}(g)$ .

*Proof.*

$$\begin{aligned}
K_{d_1, d_3}(h) &= \sup_{s_1, s_2} \frac{d_3(f(g(s_1)), f(g(s_2)))}{d_1(s_1, s_2)} \\
&= \sup_{s_1, s_2} \frac{d_2(g(s_1), g(s_2))}{d_1(s_1, s_2)} \frac{d_3(f(g(s_1)), f(g(s_2)))}{d_2(g(s_1), g(s_2))} \\
&\leq \sup_{s_1, s_2} \frac{d_2(g(s_1), g(s_2))}{d_1(s_1, s_2)} \sup_{s_1, s_2} \frac{d_3(f(s_1), f(s_2))}{d_2(s_1, s_2)} \\
&= K_{d_1, d_2}(g)K_{d_2, d_3}(f).
\end{aligned}$$

$\square$

**Lemma 3.** Given a Lipschitz function  $f : \mathcal{S} \mapsto \mathbb{R}$  with constant  $K_{d_S, d_{\mathbb{R}}}(f)$ :

$$K_{d_S, d_{\mathbb{R}}}^A \left( \int \widehat{T}(s' | s, a) f(s') ds' \right) \leq K_{d_S, d_{\mathbb{R}}}(f) K_{d_S, W}^A(\widehat{T}).$$

*Proof.*

$$\begin{aligned}
K_{d_S, d_{\mathbb{R}}}^A \left( \int_{s'} \widehat{T}(s' | s, a) f(s') ds' \right) &= \sup_a \sup_{s_1, s_2} \frac{|\int_{s'} (\widehat{T}(s' | s_1, a) - \widehat{T}(s' | s_2, a)) f(s') ds'|}{d(s_1, s_2)} \\
&= \sup_a \sup_{s_1, s_2} \frac{|\int_{s'} (\widehat{T}(s' | s_1, a) - \widehat{T}(s' | s_2, a)) f(s') \frac{K_{d_S, d_{\mathbb{R}}}(f)}{K_{d_S, d_{\mathbb{R}}}(f)} ds'|}{d(s_1, s_2)} \\
&= K_{d_S, d_{\mathbb{R}}}(f) \sup_a \sup_{s_1, s_2} \frac{|\int_{s'} (\widehat{T}(s' | s_1, a) - \widehat{T}(s' | s_2, a)) \frac{f(s')}{K_{d_S, d_{\mathbb{R}}}(f)} ds'|}{d(s_1, s_2)} \\
&\leq K_{d_S, d_{\mathbb{R}}}(f) \sup_a \sup_{s_1, s_2} \frac{|\sup_{g:K_{d_S, d_{\mathbb{R}}}(g) \leq 1} \int_{s'} (\widehat{T}(s' | s_1, a) - \widehat{T}(s' | s_2, a)) g(s') ds'|}{d(s_1, s_2)} \\
&= K_{d_S, d_{\mathbb{R}}}(f) \sup_a \sup_{s_1, s_2} \frac{\sup_{g:K_{d_S, d_{\mathbb{R}}}(g) \leq 1} \int_{s'} (\widehat{T}(s' | s_1, a) - \widehat{T}(s' | s_2, a)) g(s') ds'}{d(s_1, s_2)} \\
&= K_{d_S, d_{\mathbb{R}}}(f) \sup_a \sup_{s_1, s_2} \frac{W(\widehat{T}(\cdot | s_1, a), \widehat{T}(\cdot | s_2, a))}{d(s_1, s_2)} \\
&= K_{d_S, d_{\mathbb{R}}}(f) K_{d_S, W}^A(\widehat{T}).
\end{aligned}$$

□

**Lemma 4.** *The following operators (Asadi & Littman, 2017) are Lipschitz with constants:*

1.  $K_{\|\cdot\|_\infty, d_R}(\max(x)) = K_{\|\cdot\|_\infty, d_R}(\text{mean}(x)) = K_{\|\cdot\|_\infty, d_R}(\epsilon\text{-greedy}(x)) = 1$
2.  $K_{\|\cdot\|_\infty, d_R}(mm_\beta(x) := \frac{\log \frac{\sum_i e^{\beta x_i}}{n}}{\beta}) = 1$
3.  $K_{\|\cdot\|_\infty, d_R}(\text{boltz}_\beta(x) := \frac{\sum_{i=1}^n x_i e^{\beta x_i}}{\sum_{i=1}^n e^{\beta x_i}}) \leq \sqrt{|A|} + \beta V_{\max} |A|$

*Proof.* 1 was proven by Littman & Szepesvári (1996), and 2 is proven several times (Fox et al., 2016; Asadi & Littman, 2017; Nachum et al., 2017; Neu et al., 2017). We focus on proving 3. Define

$$\rho(x)_i = \frac{e^{\beta x_i}}{\sum_{i=1}^n e^{\beta x_i}},$$

and observe that  $\text{boltz}_\beta(x) = x^\top \rho(x)$ . Gao & Pavel (2017) showed that  $\rho$  is Lipschitz:

$$\|\rho(x_1) - \rho(x_2)\|_2 \leq \beta \|x_1 - x_2\|_2 \quad (1)$$

Using their result, we can further show:

$$\begin{aligned} & |\rho(x_1)^\top x_1 - \rho(x_2)^\top x_2| \\ & \leq |\rho(x_1)^\top x_1 - \rho(x_1)^\top x_2| + |\rho(x_1)^\top x_2 - \rho(x_2)^\top x_2| \\ & \leq \|\rho(x_1)\|_2 \|x_1 - x_2\|_2 \\ & \quad + \|x_2\|_2 \|\rho(x_1) - \rho(x_2)\|_2 \quad (\text{Cauchy-Shwartz}) \\ & \leq \|\rho(x_1)\|_2 \|x_1 - x_2\|_2 \\ & \quad + \|x_2\|_2 \beta \|x_1 - x_2\|_2 \quad (\text{from Eqn 1}) \\ & \leq (1 + \beta V_{\max} \sqrt{|A|}) \|x_1 - x_2\|_2 \\ & \leq (\sqrt{|A|} + \beta V_{\max} |A|) \|x_1 - x_2\|_\infty, \end{aligned}$$

dividing both sides by  $\|x_1 - x_2\|_\infty$  leads to 3. □

Below, we derive the Lipschitz constant for various functions.

**ReLU non-linearity** We show that  $\text{ReLU} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  has Lipschitz constant 1 for  $p$ .

$$\begin{aligned} K_{\|\cdot\|_p, \|\cdot\|_p}(\text{ReLU}) &= \sup_{x_1, x_2} \frac{\|\text{ReLU}(x_1) - \text{ReLU}(x_2)\|_p}{\|x_1 - x_2\|_p} \\ &= \sup_{x_1, x_2} \frac{(\sum_i |\text{ReLU}(x_1)_i - \text{ReLU}(x_2)_i|^p)^{\frac{1}{p}}}{\|x_1 - x_2\|_p} \\ & \quad (\text{We can show that } |\text{ReLU}(x_1)_i - \text{ReLU}(x_2)_i| \leq |x_{1,i} - x_{2,i}| \text{ and so) :} \\ &\leq \sup_{x_1, x_2} \frac{(\sum_i |x_{1,i} - x_{2,i}|^p)^{\frac{1}{p}}}{\|x_1 - x_2\|_p} \\ &= \sup_{x_1, x_2} \frac{\|x_1 - x_2\|_p}{\|x_1 - x_2\|_p} = 1 \end{aligned}$$

**Matrix multiplication** Let  $W \in \mathbb{R}^{n \times m}$ . We derive the Lipschitz continuity for the function  $\times W(x) = Wx$ .

For  $p = \infty$  we have:

$$K_{\|\cdot\|_\infty, \|\cdot\|_\infty}(\times W(x_1))$$

$$\begin{aligned}
&= \sup_{x_1, x_2} \frac{\|\times W(x_1) - \times W(x_2)\|_\infty}{\|x_1 - x_2\|_\infty} = \sup_{x_1, x_2} \frac{\|Wx_1 - Wx_2\|_\infty}{\|x_1 - x_2\|_\infty} = \sup_{x_1, x_2} \frac{\|W(x_1 - x_2)\|_\infty}{\|x_1 - x_2\|_\infty} \\
&= \sup_{x_1, x_2} \frac{\sup_j |W_j(x_1 - x_2)|}{\|x_1 - x_2\|_\infty} \\
&\leq \sup_{x_1, x_2} \frac{\sup_j \|W_j\| \|x_1 - x_2\|_\infty}{\|x_1 - x_2\|_\infty} \quad (\text{Hölder's inequality}) \\
&= \sup_j \|W_j\|_1,
\end{aligned}$$

where  $W_j$  refers to  $j$ th row of the weight matrix  $W$ . Similarly, for  $p = 1$  we have:

$$\begin{aligned}
&K_{\|\cdot\|_1, \|\cdot\|_1}(\times W(x_1)) \\
&= \sup_{x_1, x_2} \frac{\|\times W(x_1) - \times W(x_2)\|_1}{\|x_1 - x_2\|_1} = \sup_{x_1, x_2} \frac{\|Wx_1 - Wx_2\|_1}{\|x_1 - x_2\|_1} = \sup_{x_1, x_2} \frac{\|W(x_1 - x_2)\|_1}{\|x_1 - x_2\|_1} \\
&= \sup_{x_1, x_2} \frac{\sum_j |W_j(x_1 - x_2)|}{\|x_1 - x_2\|_1} \\
&\leq \sup_{x_1, x_2} \frac{\sum_j \|W_j\|_\infty \|x_1 - x_2\|_1}{\|x_1 - x_2\|_1} = \sum_j \|W_j\|_\infty,
\end{aligned}$$

and finally for  $p = 2$ :

$$\begin{aligned}
&K_{\|\cdot\|_2, \|\cdot\|_2}(\times W(x_1)) \\
&= \sup_{x_1, x_2} \frac{\|\times W(x_1) - \times W(x_2)\|_2}{\|x_1 - x_2\|_2} = \sup_{x_1, x_2} \frac{\|Wx_1 - Wx_2\|_2}{\|x_1 - x_2\|_2} = \sup_{x_1, x_2} \frac{\|W(x_1 - x_2)\|_2}{\|x_1 - x_2\|_2} \\
&= \sup_{x_1, x_2} \frac{\sqrt{\sum_j |W_j(x_1 - x_2)|^2}}{\|x_1 - x_2\|_2} \\
&\leq \sup_{x_1, x_2} \frac{\sqrt{\sum_j \|W_j\|_2^2 \|x_1 - x_2\|_2^2}}{\|x_1 - x_2\|_2} = \sqrt{\sum_j \|W_j\|_2^2}.
\end{aligned}$$

**Vector addition** We show that  $+b : \mathbb{R}^n \rightarrow \mathbb{R}^n$  has Lipschitz constant 1 for  $p = 0, 1, \infty$  for all  $b \in \mathbb{R}^n$ .

$$\begin{aligned}
K_{\|\cdot\|_p, \|\cdot\|_p}(\text{ReLU}) &= \sup_{x_1, x_2} \frac{\|+b(x_1) - +b(x_2)\|_p}{\|x_1 - x_2\|_p} \\
&= \sup_{x_1, x_2} \frac{\|(x_1 + b) - (x_2 + b)\|_p}{\|x_1 - x_2\|_p} = \frac{\|x_1 - x_2\|_p}{\|x_1 - x_2\|_p} = 1
\end{aligned}$$

**Supervised-learning domain** We used the following 5 functions to generate the dataset:

$$\begin{aligned}
f_0(x) &= \tanh(x) + 3 \\
f_1(x) &= x * x \\
f_2(x) &= \sin(x) - 5 \\
f_3(x) &= \sin(x) - 3 \\
f_4(x) &= \sin(x) * \sin(x)
\end{aligned}$$

We sampled each function 30 times, where the input was chosen uniformly randomly from  $[-2, 2]$  each time.

---

## References

- Asadi, K. and Littman, M. L. An alternative softmax operator for reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 243–252, 2017.
- Fox, R., Pakman, A., and Tishby, N. G-learning: Taming the noise in reinforcement learning via soft updates. *Uncertainty in Artificial Intelligence*, 2016.
- Gao, B. and Pavel, L. On the properties of the softmax function with application in game theory and reinforcement learning. *arXiv preprint arXiv:1704.00805*, 2017.
- Littman, M. L. and Szepesvári, C. A generalized reinforcement-learning model: Convergence and applications. In *Proceedings of the 13th International Conference on Machine Learning*, pp. 310–318, 1996.
- Nachum, O., Norouzi, M., Xu, K., and Schuurmans, D. Bridging the gap between value and policy based reinforcement learning. *arXiv preprint arXiv:1702.08892*, 2017.
- Neu, G., Jonsson, A., and Gómez, V. A unified view of entropy-regularized Markov decision processes. *arXiv preprint arXiv:1705.07798*, 2017.