# On-line Reinforcement Learning Using Incremental Kernel-Based Stochastic Factorization

## SUPPLEMENTARY MATERIAL

**André M. S. Barreto**
School of Computer Science
McGill University
Montreal, Canada
amsb@cs.mcgill.ca

**Doina Precup**
School of Computer Science
McGill University
Montreal, Canada
dprecup@cs.mcgill.ca

**Joelle Pineau**
School of Computer Science
McGill University
Montreal, Canada
jpineau@cs.mcgill.ca

### Abstract

This is the supplementary material for the paper entitled "On-line Reinforcement Learning Using Incremental Kernel-Based Stochastic Factorization" [2]. It contains the details of our theoretical developments that could not be included in the paper due to space constraints. This material should be read in conjunction with the main paper.

## 1 Preliminaries

- Similarly to Ormoneit and Sen [3], we define a "mother kernel" $\phi(x) : \mathbb{R}^+ \mapsto \mathbb{R}^+$ satisfying

  (i)  $\phi(x)$ is continuous in $\mathbb{R}^+$,

  (ii)  $\int_0^\infty \phi(x)dx \le L_\phi < \infty$,

  (iii)  $\phi(x) \ge \phi(y)$ if $x < y$,

  (iv)  $\exists A_\phi, \lambda_\phi > 0, \exists B_\phi \ge 0$ such that $A_\phi \exp(-x) \le \phi(x) \le \lambda_\phi A_\phi \exp(-x)$ if $x \ge B_\phi$.

  **Remarks:**

    - Assumption (i) is implied by Ormoneit and Sen's [3] assumption that $\phi$ is Lipschitz continuous. Ormoneit and Sen also assume that $\int_0^1 \phi(z)dz = 1$ (see Appendix A.1 in [3]).
    - Assumption (iv) implies that the kernel function $\phi$ will eventually decay exponentially and also that $\phi(z) > 0$ for all $z \in \mathbb{R}^+$.

- Let $\mathbb{S} \subset [0,1]^d$ and let $\| \cdot \|$ be a norm in $\mathbb{R}^d$. Then, we define

$$k_\tau(s, s') = \phi\left(\frac{\| s - s' \|}{\tau}\right),$$

  where $\tau > 0$ is the "width" of the kernel $k_\tau$.

- Let $M$ be a Markov decision process (MDP) with state space $\mathbb{S}$ and let $S^a = \{(s_k^a, r_k^a, \hat{s}_k^a)|k = 1, 2, ..., n_a\}$ be a set of sample transitions associated with action $a \in A$, where $s_k^a, \hat{s}_k^a \in \mathbb{S}$ and $r_k^a \in \mathbb{R}$. We define the normalized kernel function associated with action $a$ as

$$\kappa_\tau^a(s, s_i^a) = \frac{k_\tau(s, s_i^a)}{\sum_{j=1}^{n_a} k_\tau(s, s_j^a)}.$$

- Let $\bar{S} \equiv \{\bar{s}_1, \bar{s}_2, ..., \bar{s}_m\}$ be a set of representative states in $\mathbb{S}$. Define:
    - $\hat{s}_*^a \equiv \hat{s}_k^a$ with $k = \operatorname{argmax}_i \min_j \| \hat{s}_i^a - \bar{s}_j \|$,

- $\bar{s}^a_* \equiv \bar{s}_h$ where $h = \text{argmin}_j \parallel \hat{s}^a_* - \bar{s}_j \parallel$,
- $\hat{s}_* \equiv \hat{s}^b_*$ where $b = \text{argmax}_a \parallel \hat{s}^a_* - \bar{s}^a_* \parallel$,
- $\bar{s}_* \equiv \bar{s}^b_*$ where $b = \text{argmax}_a \parallel \hat{s}^a_* - \bar{s}^a_* \parallel$,
- $\mathfrak{d}^* \equiv \parallel \hat{s}_* - \bar{s}_* \parallel$.

We assume that

(v) $\hat{s}^a_*$ and $\bar{s}^a_*$ are unique for all $a \in A$.

# 2 Data-independent definitions

**Definition 1.** *For any $\alpha \in (0,1]$, the $\alpha$-radius of $k_\tau$ with respect to $s$ and $s'$ is defined as*

$$\rho(k_\tau, s, s', \alpha) = \max\left\{x \in \mathbb{R}^+ | \phi\left(\frac{x}{\tau}\right) = \alpha k_\tau(s, s')\right\}.$$

**Remarks:**

- The existence of $\rho(k_\tau, s, s', \alpha)$ is guaranteed by properties (i), (ii) and (iii).
- $\rho(k_\tau, s, s', \alpha) \geq \parallel s - s' \parallel$.

**Property 1.** *If $\parallel s - s' \parallel < \parallel s - s'' \parallel$, then $\rho(k_\tau, s, s', \alpha) \leq \rho(k_\tau, s, s'', \alpha)$.*

*Proof.* Let $r = \rho(k_\tau, s, s', \alpha)$. Then,

$$\phi\left(\frac{r}{\tau}\right) = \alpha k_\tau(s, s') = \alpha\phi\left(\frac{\parallel s - s' \parallel}{\tau}\right) \geq \alpha\phi\left(\frac{\parallel s - s'' \parallel}{\tau}\right) = \alpha k_\tau(s, s'').$$

If $\phi(r/\tau) = \alpha k_\tau(s, s'')$, then $\rho(k_\tau, s, s'', \alpha) = r$. If $\phi(r/\tau) > \alpha k_\tau(s, s'')$, then from (iii) it must be the case that $r < \rho(k_\tau, s, s'', \alpha)$. $\qquad\square$

**Property 2.** *If $\alpha < \alpha'$, then $\rho(k_\tau, s, s', \alpha) > \rho(k_\tau, s, s', \alpha')$.*

*Proof.* Let $r = \rho(k_\tau, s, s', \alpha')$. Then,

$$\phi\left(\frac{r}{\tau}\right) = \alpha' k_\tau(s, s') = \alpha'\phi\left(\frac{\parallel s - s' \parallel}{\tau}\right) > \alpha\phi\left(\frac{\parallel s - s' \parallel}{\tau}\right) = \alpha k_\tau(s, s').$$

From (iii) it must be the case that $r < \rho(k_\tau, s, s', \alpha)$. $\qquad\square$

**Property 3.** *For any $\alpha \in (0,1)$ and any $\varepsilon > 0$, there is a $\delta > 0$ such that $\rho(k_\tau, s, s', \alpha) - \parallel s - s' \parallel < \varepsilon$ if $\tau < \delta$.*

*Proof.* Let $z = \parallel s - s' \parallel$. We will show that, for any $\varepsilon > 0$, there is a $\delta > 0$ such that $\phi((z+\varepsilon)/\tau) < \alpha\phi(z/\tau)$ if $\tau < \delta$. We know that

$$\frac{\exp(-(z+\varepsilon)/\tau)}{\exp(-z/\tau)} < \alpha/\lambda_\phi \iff \ln\left(\frac{\exp(-(z+\varepsilon)/\tau)}{\exp(-z/\tau)}\right) < \ln(\alpha/\lambda_\phi) \iff$$

$$\iff -\frac{\varepsilon}{\tau} < \ln(\alpha/\lambda_\phi) \iff \tau < -\frac{\varepsilon}{\ln(\alpha/\lambda_\phi)}$$

(note that it must be the case that $\alpha/\lambda_\phi \neq 1$). Thus, by taking $\delta < \min(-\varepsilon/\ln(\alpha/\lambda_\phi), z/B_\phi)$ and resorting to Assumption (iv), we can write:

$$\begin{aligned}
\alpha/\lambda_\phi &> \frac{\exp(-(z+\varepsilon)/\delta)}{\exp(-z/\delta)} = \frac{A_\phi \exp(-(z+\varepsilon)/\delta)}{A_\phi \exp(-z/\delta)} \\
&\geq \frac{A_\phi \exp(-(z+\varepsilon)/\delta)}{\phi(z/\delta)} = \frac{\lambda_\phi A_\phi \exp(-(z+\varepsilon)/\delta)}{\lambda_\phi \phi(z/\delta)} \\
&\geq \frac{\phi((z+\varepsilon)/\delta)}{\lambda_\phi \phi(z/\delta)},
\end{aligned}$$

and therefore $\dfrac{\phi((z+\varepsilon)/\tau)}{\phi(z/\tau)} < \alpha$ if $\tau \leq \delta$. $\qquad\square$

**Remarks:**

- $\rho(k_\tau, s, s', \alpha) - \| s - s' \| < \varepsilon$ if $\tau < \min(-\varepsilon/\ln(\alpha/\lambda_\phi), \| s - s' \| / B_\phi)$, where $\lambda_\phi$ and $B_\phi$ depend on the particular choice of function $\phi$ (see Assumption (iv)).

- Given $s, s'$, and $s''$, with $\| s - s' \| < \| s - s'' \|$, Property 3 states that for any $\alpha \in (0, 1)$, there is a $\delta > 0$ such that $k_\tau(s, s'') < \alpha k_\tau(s, s')$ if $\tau < \delta$ (to see why this is so, it suffices to make $\varepsilon = \| s - s'' \| - \| s - s' \|$).

# 3  Data-dependent definitions

**Definition 2.** *Given $\beta > 0$, the $\beta$-dissimilarity between $s$ and $s'$ with respect to $\kappa_\tau^a$ is defined as*

$$\psi(\kappa_\tau^a, s, s', \beta) = \begin{cases} \sum_{k=1}^{n_a} |\kappa_\tau^a(s, s_k^a) - \kappa_\tau^a(s', s_k^a)|, & \text{if } \| s - s' \| \le \beta, \\ 0, & \text{otherwise}. \end{cases}$$

**Remark:** $\psi(\kappa_\tau^a, s, s', \beta) \in [0, 2]$.

**Property 4.** *For any $\beta > 0$ and any $\varepsilon > 0$, there is a $\delta > 0$ such that $\psi(\kappa_\tau^a, s, s', \beta) < \varepsilon$ if $\| s - s' \| < \delta$.*

*Proof.* If $\beta < \| s - s' \|$, then $\psi(\kappa_\tau^a, s, s', \beta) = 0$ and the result follows (see Definition 2). Otherwise:

$$\psi(\kappa_\tau^a, s, s', \beta) \equiv \psi_{\tau, s, \beta}^a(s') = \sum_{k=1}^{n_a} \left| \frac{k_\tau(s, s_k^a)}{\sum_{l=1}^{n_a} k_\tau(s, s_l^a)} - \frac{k_\tau(s', s_k^a)}{\sum_{l=1}^{n_a} k_\tau(s', s_l^a)} \right|$$

$$= \sum_{k=1}^{n_a} \left| \frac{\phi\left(\| s - s_k^a \| / \tau\right)}{\sum_{l=1}^{n_a} \phi\left(\| s - s_l^a \| / \tau\right)} - \frac{\phi\left(\| s' - s_k^a \| / \tau\right)}{\sum_{l=1}^{n_a} \phi\left(\| s' - s_l^a \| / \tau\right)} \right|.$$

From the definition of $\phi$, it is obvious that $\psi_{\tau, s, \beta}^a(s')$ is continuous in $s'$. The property follows from the fact that $\lim_{s' \to s} \psi_{\tau, s, \beta}^a(s') = 0$. $\qquad\square$

**Remarks:**

- $\psi(\kappa_\tau^a, s, s', \beta)$ does *not* necessarily increase with $\| s - s' \|$.

- Given $\varepsilon > 0$, $\delta$ is data-dependent.

# 4  Main Results

**Lemma 1.** *For any $\alpha \in (0, 1]$ and any $t \ge m - 1$, let $\delta_a = \rho(k_{\bar{\tau}}, \hat{s}_*^a, \bar{s}_*^a, \alpha/t)$, let $\psi_\delta^a = \max_{i,j} \psi(\kappa_\tau^a, \hat{s}_i^a, \bar{s}_j, \delta_a)$ and let $\psi_{\max}^a = \max_{i,j} \psi(\kappa_\tau^a, \hat{s}_i^a, \bar{s}_j, \infty)$. Then,*

$$\| \mathbf{P}^a - \mathbf{D K}^a \|_\infty \le \frac{1}{1 + \alpha} \psi_\delta^a + \frac{\alpha}{1 + \alpha} \psi_{\max}^a. \tag{1}$$

*Proof.* Let $\tilde{\mathbf{P}}^a = \mathbf{D}\mathbf{K}^a$. Recalling that $\kappa_\tau^a(s, s_j^a) = \dfrac{k_\tau(s, s_j^a)}{\sum_{k=1}^{n_a} k_\tau(s, s_k^a)}$, we can write:

$$
\begin{aligned}
\|\mathbf{p}_i^a - \tilde{\mathbf{p}}_i^a\|_1 &= \sum_{j=1}^{n_a} \left| p_{ij}^a - \sum_{k=1}^m \dot{d}_{ik}^a \dot{k}_{kj}^a \right| \\
&= \sum_{j=1}^{n_a} \left| \frac{k_\tau(\hat{s}_i^a, s_j^a)}{\sum_{l=1}^n k_\tau(\hat{s}_i^a, s_l^a)} - \sum_{k=1}^m \left( \frac{k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k)}{\sum_{l=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \frac{k_\tau(\bar{s}_k, s_j^a)}{\sum_{l=1}^n k_\tau(\bar{s}_k, s_l^a)} \right) \right| \\
&= \sum_{j=1}^{n_a} \left| \frac{k_\tau(\hat{s}_i^a, s_j^a)}{\sum_{l=1}^n k_\tau(\hat{s}_i^a, s_l^a)} - \frac{1}{\sum_{l=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \sum_{k=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k) \frac{k_\tau(\bar{s}_k, s_j^a)}{\sum_{l=1}^n k_\tau(\bar{s}_k, s_l^a)} \right| \\
&= \sum_{j=1}^{n_a} \left| \frac{1}{\sum_{l=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \sum_{k=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k) \frac{k_\tau(\hat{s}_i^a, s_j^a)}{\sum_{l=1}^n k_\tau(\hat{s}_i^a, s_l^a)} - \frac{1}{\sum_{l=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \sum_{k=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k) \frac{k_\tau(\bar{s}_k, s_j^a)}{\sum_{l=1}^n k_\tau(\bar{s}_k, s_l^a)} \right| \\
&\leq \sum_{j=1}^{n_a} \frac{1}{\sum_{l=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \sum_{k=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k) \left| \frac{k_\tau(\hat{s}_i^a, s_j^a)}{\sum_{l=1}^n k_\tau(\hat{s}_i^a, s_l^a)} - \frac{k_\tau(\bar{s}_k, s_j^a)}{\sum_{l=1}^n k_\tau(\bar{s}_k, s_l^a)} \right| \\
&= \frac{1}{\sum_{l=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \sum_{k=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k) \sum_{j=1}^{n_a} \left| \kappa_\tau^a(\hat{s}_i^a, s_j^a) - \kappa_\tau^a(\bar{s}_k, s_j^a) \right|.
\end{aligned}
$$

Let $H = \{k \mid \| \hat{s}_i^a - \bar{s}_k \| \leq \delta_a\}$ and let $\bar{H} = \{1, 2, ..., m\} - H$. Then,

$$
\begin{aligned}
\|\mathbf{p}_i^a - \tilde{\mathbf{p}}_i^a\|_1 &\leq \sum_{k \in H} \frac{k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k)}{\sum_{l=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \sum_{j=1}^{n_a} \left| \kappa_\tau^a(\hat{s}_i^a, s_j^a) - \kappa_\tau^a(\bar{s}_k, s_j^a) \right| + \sum_{k \in \bar{H}} \frac{k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k)}{\sum_{l=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \sum_{j=1}^{n_a} \left| \kappa_\tau^a(\hat{s}_i^a, s_j^a) - \kappa_\tau^a(\bar{s}_k, s_j^a) \right| \\
&\leq \frac{\sum_{k \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k)}{\sum_{l=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \max_{h \in H} \psi(\kappa_\tau^a, \hat{s}_i^a, \bar{s}_h, \infty) + \frac{\sum_{k \in \bar{H}} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k)}{\sum_{l=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \max_{h \in \bar{H}} \psi(\kappa_\tau^a, \hat{s}_i^a, \bar{s}_h, \infty) \\
&\leq \frac{\sum_{k \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k)}{\sum_{l=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \psi_\delta^a + \frac{\sum_{k \in \bar{H}} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k)}{\sum_{l=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \max_h \psi(\kappa_\tau^a, \hat{s}_i^a, \bar{s}_h, \infty) \\
&\leq \frac{\sum_{k \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k)}{\sum_{l=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \psi_\delta^a + \frac{\sum_{k \in \bar{H}} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k)}{\sum_{l=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \psi_{\max}^a \\
&= \frac{\sum_{k \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k)}{\sum_{l \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l) + \sum_{l \in \bar{H}} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \psi_\delta^a + \frac{\sum_{k \in \bar{H}} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k)}{\sum_{l \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l) + \sum_{l \in \bar{H}} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l)} \psi_{\max}^a.
\end{aligned}
\tag{2}
$$

From the definition of $\delta_a$ we can write

$$
k_{\bar{\tau}}(\hat{s}_*^a, \bar{s}_k) \leq \frac{\alpha}{t} k_{\bar{\tau}}(\hat{s}_*^a, \bar{s}_*^a) \text{ if } \| \hat{s}_*^a - \bar{s}_k \| \geq \delta_a. \tag{3}
$$

Now, let $w = \arg\min_k \| \hat{s}_i^a - \bar{s}_k \|$. We know that $\| \hat{s}_i^a - \bar{s}_w \| \leq \| \hat{s}_*^a - \bar{s}_*^a \|$, and thus, from Assumption (iii), it follows that $k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w) \geq k_{\bar{\tau}}(\hat{s}_*^a, \bar{s}_*^a)$. This fact together with (3) imply that

$$
k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k) \leq \frac{\alpha}{t} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w) \text{ if } \bar{s}_w \in \bar{H},
$$

which allows us to write

$$
\sum_{l \in \bar{H}} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l) \leq \frac{\alpha |\bar{H}|}{t} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w). \tag{4}
$$

Plugging (4) back into (2), we can write:

$$
\begin{aligned}
\|\mathbf{p}_i^a - \tilde{\mathbf{p}}_i^a\|_1 &\leq \frac{\sum_{k \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k)}{\sum_{l \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l) + \frac{\alpha |\bar{H}|}{t} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w)} \psi_\delta^a + \frac{\frac{\alpha |\bar{H}|}{t} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w)}{\sum_{l \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l) + \frac{\alpha |\bar{H}|}{t} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w)} \psi_{\max}^a \tag{5} \\
&\leq \frac{\sum_{k \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k)}{\sum_{l \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l) + \frac{\alpha t}{t} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w)} \psi_\delta^a + \frac{\frac{\alpha t}{t} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w)}{\sum_{l \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l) + \frac{\alpha t}{t} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w)} \psi_{\max}^a \tag{6} \\
&= \frac{\sum_{k \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k)}{\sum_{l \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l) + \alpha k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w)} \psi_\delta^a + \frac{\alpha k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w)}{\sum_{l \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l) + \alpha k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w)} \psi_{\max}^a, \tag{7}
\end{aligned}
$$

where in (5) and (6) we used the fact that the coefficients multiplying $\psi_\delta^a$ and $\psi_{\max}^a$ define a convex combination, and we are increasing the weight of the latter. Noticing that

$$\frac{\alpha k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w)}{\sum_{l \in H} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_l) + \alpha k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w)} \leq \frac{\alpha k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w)}{k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w) + \alpha k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w)} = \frac{\alpha}{1+\alpha},$$

and applying the same reasoning to the coefficients of $\psi_\delta^a$ and $\psi_{\max}^a$ in (7), we can finally write

$$\|\mathbf{p}_i^a - \tilde{\mathbf{p}}_i^a\|_1 \leq \frac{1}{1+\alpha} \psi_\delta^a + \frac{\alpha}{1+\alpha} \psi_{\max}^a.$$

$\square$

**Remarks:**

- $\psi_\delta^a \to \psi_{\max}^a$ as $\alpha \to 0$.
- There is an $\alpha^* \in (0, 1]$ that minimizes the right-hand side of (1).

Let $\tilde{\mathbf{v}} = \Gamma \mathbf{D} \bar{\mathbf{Q}}^*$, where $\Gamma$ is the 'max' operator applied row wise, that is, $\tilde{v}_i = \max_a (\mathbf{D} \bar{\mathbf{Q}}^*)_{ia}$. Recalling that $\mathfrak{d}^*$ is the maximum distance from a sampled state $\hat{s}_i^a$ to the closest representative state and that $\bar{\tau}$ is the width of kernel $k_{\bar{\tau}}$, we present the following result:

**Proposition 1.** *For any $\varepsilon > 0$, there are $\delta_1, \delta_2 > 0$ such that $\|\hat{\mathbf{v}}^* - \tilde{\mathbf{v}}\|_\infty < \varepsilon$ if $\mathfrak{d}^* < \delta_1$ and $\bar{\tau} < \delta_2$.*

*Proof.* We have previously showed that

$$\|\hat{\mathbf{v}}^* - \tilde{\mathbf{v}}\|_\infty \leq \frac{1}{1-\gamma} \max_a \|\hat{\mathbf{r}}^a - \mathbf{D}\bar{\mathbf{r}}^a\|_\infty + \frac{1}{(1-\gamma)^2} \left( \bar{C} \max_i (1 - \max_j d_{ij}) + \frac{\hat{C}\gamma}{2} \max_a \|\hat{\mathbf{P}}^a - \mathbf{D}\mathbf{K}^a\|_\infty \right), \quad (8)$$

where $\|\cdot\|_\infty$ is the infinity norm, $\hat{\mathbf{v}}^* \in \mathbb{R}^n$ is the optimal value function of KBRL's MDP, $\hat{C} = \max_{a,i} \hat{r}_i^a - \min_{a,i} \hat{r}_i^a$, $\bar{C} = \max_{a,i} \bar{r}_i^a - \min_{a,i} \bar{r}_i^a$, and $\mathbf{K}^a$ is matrix $\mathbf{K}$ with all elements equal to zero except for those corresponding to matrix $\dot{\mathbf{K}}^a$ (see [1, 2] for details). Let $\check{\mathbf{r}} \equiv [(\mathbf{r}^1)^\top, (\mathbf{r}^2)^\top, ..., (\mathbf{r}^{|A|})^\top]^\top \in \mathbb{R}^n$, where $\mathbf{r}^a \in \mathbb{R}^{n_a}$ is the vector composed of sample rewards $r_i^a$. Then,

$$\|\hat{\mathbf{r}}^a - \mathbf{D}\bar{\mathbf{r}}^a\|_\infty = \|\hat{\mathbf{P}}^a \check{\mathbf{r}} - \mathbf{D}\dot{\mathbf{K}}^a \mathbf{r}^a\|_\infty = \|\hat{\mathbf{P}}^a \check{\mathbf{r}} - \mathbf{D}\mathbf{K}^a \check{\mathbf{r}}\|_\infty = \|(\hat{\mathbf{P}}^a - \mathbf{D}\mathbf{K}^a)\check{\mathbf{r}}\|_\infty \leq \|\hat{\mathbf{P}}^a - \mathbf{D}\mathbf{K}^a\|_\infty \|\check{\mathbf{r}}\|_\infty, \quad (9)$$

where the equality $\hat{\mathbf{r}}^a = \hat{\mathbf{P}}^a \check{\mathbf{r}}$ is a consequence of the fact that KBRL's reward function $R^a(s, s')$ is independent of the start state $s$ (see (1) in the main paper [2]). Thus, plugging (9) back into (8), it is clear that there is a $\eta > 0$ such that $\|\hat{\mathbf{v}}^* - \tilde{\mathbf{v}}\|_\infty < \varepsilon$ if $\max_a \|\hat{\mathbf{P}}^a - \mathbf{D}\mathbf{K}^a\|_\infty < \eta$ and $\max_i (1 - \max_j d_{ij}) < \eta$. We start by showing that if $\mathfrak{d}^*$ and $\bar{\tau}$ are small enough, then $\max_a \|\hat{\mathbf{P}}^a - \mathbf{D}\mathbf{K}^a\|_\infty < \eta$. From Lemma 1 we know that, for any set of $m \leq n$ representative states, and for any $\alpha \in (0, 1]$, the following must hold:

$$\max_a \|\mathbf{P}^a - \mathbf{D}\mathbf{K}^a\|_\infty \leq \frac{1}{1+\alpha} \psi_\rho + \frac{\alpha}{1+\alpha} \psi_{\text{MAX}},$$

where $\psi_{\text{MAX}} = \max_{a,i,s} \psi(k_\tau, \hat{s}_i^a, s, \infty)$ and $\psi_\rho = \max_a \psi_\rho^a = \max_{a,i,j} \psi(\kappa_\tau^a, \hat{s}_i^a, \bar{s}_j, \rho^a)$, with $\rho^a = \rho(k_{\bar{\tau}}, \hat{s}_*^a, \bar{s}_*^a, \alpha/(n-1))$. Note that $\psi_{\text{MAX}}$ is independent of the representative states. Define $\alpha$ such that $\alpha/(1+\alpha)\psi_{\text{MAX}} < \eta$. We have to show that, if we define the representative states in such a way that $\mathfrak{d}^*$ is small enough, and set $\bar{\tau}$ accordingly, then we can make $\psi_\rho < (1-\alpha)\eta - \alpha\psi_{\text{MAX}} \equiv \eta'$. From Property 4 we know that there is a $\delta_1 > 0$ such that $\psi_\rho < \eta'$ if $\rho^a < \delta_1$ for all $a \in A$. From Property 1 we know that $\rho^a \leq \rho(k_{\bar{\tau}}, \hat{s}_*, \bar{s}_*, \alpha/(n-1))$ for all $a \in A$. From Property 3 we know that, for any $\varepsilon' > 0$, there is a $\delta' > 0$ such that $\rho(k_{\bar{\tau}}, \hat{s}_*, \bar{s}_*, \alpha/(n-1)) < \mathfrak{d}^* + \varepsilon'$ if $\bar{\tau} < \delta'$. Therefore, if $\mathfrak{d}^* < \delta_1$, we can take any $\varepsilon' < \delta_1 - \mathfrak{d}^*$ to have an upper bound $\delta'$ for $\bar{\tau}$. It remains to show that there is a $\delta > 0$ such that $\min_i \max_j d_{ij} > 1 - \eta$ if $\bar{\tau} < \delta$. Recalling that $d_{ij}^a = k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_j)/\sum_{k=1}^m k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_k)$, let $w = \arg\max_j k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_j)$, and let $y_i^a = k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_w)$ and $\breve{y}_i^a = \max_{j \neq w} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_j)$. Then, for any $i$,

$$\max_j d_{ij}^a = \frac{y_i^a}{\left(y_i^a + \sum_{j \neq w} k_{\bar{\tau}}(\hat{s}_i^a, \bar{s}_j)\right)} \geq \frac{y_i^a}{(y_i^a + (m-1)\breve{y}_i^a)}.$$

From Assumption (v) and Property 3 we know that there is a $\delta_i^a > 0$ such that $y_i^a > (m-1)(1-\eta)\breve{y}_i^a/\eta$ if $\bar{\tau} < \delta_i^a$. Thus, by making $\delta = \min_{a,i} \delta_i^a$, we can guarantee that $\min_i \max_j d_{ij} > 1 - \eta$. Finally, if we take $\delta_2 = \min(\delta, \delta')$, the result follows. $\square$

**Remark:** If we define a "net" over $\mathbb{S}$ using the representative states, then we know that $\mathfrak{d}^*$ is smaller than the resolution of the net.

# References

[1] A. M. S. Barreto, D. Precup, and J. Pineau. Reinforcement learning using kernel-based stochastic factorization. In *Advances in Neural Information Processing Systems (NIPS)*, pages 720–728, 2011.

[2] A. M. S. Barreto, D. Precup, and J. Pineau. On-line reinforcement learning using incremental kernel-based stochastic factorization. In *Advances in Neural Information Processing Systems (NIPS)*, 2012.

[3] D. Ormoneit and S. Sen. Kernel-based reinforcement learning. *Machine Learning*, 49 (2–3): 161–178, 2002.