# A Kernel for Hierarchical Parameter Spaces

Frank Hutter and Michael A. Osborne

`fh@informatik.uni-freiburg.de` and `mosb@robots.ox.ac.uk`

September 6, 2013

**Abstract**

We define kernels for mixed continuous/discrete spaces and conditional spaces and show that they are positive definite.

## 1 Introduction

We aim to do inference about some function $g$ with domain (input space) $\mathcal{X}$. $\mathcal{X} = \prod_{i=1}^{D} \mathcal{X}_i$ is a $D$-dimensional input space, where each individual dimension is either bounded real or categorical, that is, $\mathcal{X}_i$ is either $[l_i, u_i] \subset \mathbb{R}$ (with lower and upper bounds $l_i$ and $u_i$, respectively) or $\{v_{i,1}, \ldots, v_{i,m_i}\}$.

Associated with $\mathcal{X}$, there is a DAG structure $\mathcal{D}$, whose vertices are the dimensions $\{1, \ldots, D\}$. $\mathcal{X}$ will be restricted by $\mathcal{D}$: if vertex $i$ has children under $\mathcal{D}$, $\mathcal{X}_i$ must be categorical. $\mathcal{D}$ is also used to specify when each input is *active* (that is, relevant to inference about $g$). In particular, we assume each input dimension is only active under some instantiations of its ancestor dimensions in $\mathcal{D}$. More precisely, we define $D$ functions $\delta_i \colon \mathcal{X} \to \mathcal{B}$, for $i \in \{1, \ldots, D\}$, and where $\mathcal{B} = \{\text{true}, \text{false}\}$. We take

$$\delta_i(\underline{x}) = \delta_i\big(\underline{x}(\text{anc}_i)\big), \tag{1}$$

where $\text{anc}_i$ are the ancestor vertices of $i$ in $\mathcal{D}$, such that $\delta_i(\underline{x})$ is true only for appropriate values of those entries of $\underline{x}$ corresponding to ancestors of $i$ in $\mathcal{D}$. We say $i$ is active for $\underline{x}$ iff $\delta_i(\underline{x})$.

Our aim is to specify a kernel for $\mathcal{X}$, *i.e.*, a positive semi-definite function $k \colon \mathcal{X} \times \mathcal{X} \to \mathbb{R}$. We will first specify an individual kernel for each input dimension, *i.e.*, a positive semi-definite function $k_i \colon \mathcal{X} \times \mathcal{X} \to \mathbb{R}$. $k$ can then be taken as either a sum,

$$k(\underline{x}, \underline{x}') = \sum_{i=1}^{D} k_i(\underline{x}, \underline{x}'), \tag{2}$$

product,

$$k(\underline{x}, \underline{x}') = \prod_{i=1}^{D} k_i(\underline{x}, \underline{x}'), \tag{3}$$

or any other permitted combination, of these individual kernels. Note that each individual kernel $k_i$ will depend on an input vector $\underline{x}$ only through dependence on $x_i$ and $\delta_i(\underline{x})$,

$$k_i(\underline{x}, \underline{x}') = \tilde{k}_i\big(x_i, \delta_i(\underline{x}), x_i', \delta_i(\underline{x}')\big). \tag{4}$$

That is, $x_j$ for $j \neq i$ will influence $k_i(\underline{x}, \underline{x}')$ only if $j \in \mathrm{anc}_i$, and only by affecting whether $i$ is active.

Below we will construct pseudometrics $d_i \colon \mathcal{X} \times \mathcal{X} \to \mathbb{R}^+$: that is, $d_i$ satisfies the requirements of a metric aside from the identity of indiscernibles. As for $k_i$, these pseudometrics will depend on an input vector $\underline{x}$ only through dependence on both $x_i$ and $\delta_i(\underline{x})$. $d_i(\underline{x}, \underline{x}')$ will be designed to provide an intuitive measure of how different $g(\underline{x})$ is from $g(\underline{x}')$. For each $i$, we will then construct a (pseudo-)isometry $f_i$ from $\mathcal{X}$ to a Euclidean space ($\mathbb{R}^2$ for bounded real parameters, and $\mathbb{R}^m$ for categorical-valued parameters with $m$ choices). That is, denoting the Euclidean metric on the appropriate space as $d_E$, $f_i$ will be such that

$$d_i(\underline{x}, \underline{x}') = d_\mathrm{E}(f_i(\underline{x}), f_i(\underline{x}')) \tag{5}$$

for all $\underline{x}, \underline{x}' \in \mathcal{X}$. We can then use our transformed inputs, $f_i(\underline{x})$, within any standard Euclidean kernel $\kappa$. We'll make this explicit in Proposition 2.

---

FH: I tried to make things a bit more formal here by using explicit definitions and referring back to them later. We might want to go one step further and do something similar for isometries/pseudometrics, once it is clear that we actually need these concepts in the formal proofs. What do you think?

---

MO: sure

---

.

**Definition 1.** *A function* $\kappa \colon \mathbb{R}^+ \to \mathbb{R}$ *is* a positive semi-definite covariance function over Euclidean space *if* $K \in \mathbb{R}^{N \times N}$*, defined by*

$$K_{m,n} = \kappa\big(d_E(\underline{y}_m, \underline{y}_n)\big), \quad \textit{for } \underline{y}_m, \underline{y}_n \in \mathbb{R}^P, \quad m, n = 1, \dots, N,$$

*is positive semi-definite for any* $\underline{y}_1, \dots, \underline{y}_N \in \mathbb{R}^P$.

A popular example of such a $\kappa$ is the exponentiated quadratic, for which $\kappa(\delta) = \sigma^2 \exp(-\frac{1}{2} \frac{\delta^2}{\lambda^2})$; another popular choice is the rational quadratic, for which $\kappa(\delta) = \sigma^2 (1 + \frac{1}{2\alpha} \frac{\delta^2}{\lambda^2})^{-\alpha}$.

**Proposition 2.** *Let* $\kappa$ *be a positive semi-definite covariance function over Euclidean space and let* $d_i$ *satisfy Equation 5. Then,* $k_i \colon \mathcal{X} \times \mathcal{X} \to \mathbb{R}^+$*, defined by*

$$k_i(\underline{x}, \underline{x}') = \kappa\big(d_i(\underline{x}, \underline{x}')\big)$$

*is a positive semi-definite covariance function over input space* $\mathcal{X}$.

*Proof.* We need to show that for any $\underline{x}_1, \ldots, \underline{x}_N \in \mathcal{X}$, $K \in \mathbb{R}^{N \times N}$ defined by

$$K_{m,n} = \kappa\big(d_i(\underline{x}_m, \underline{x}_n)\big), \quad \text{for } \underline{x}_m, \underline{x}_n \in \mathcal{X}, \quad m, n = 1, \ldots, N,$$

is positive semi-definite. Now, by the definition of $d_i$,

$$K_{m,n} = \kappa\Big(d_{\mathrm{E}}(f_i(\underline{x}_m), f_i(\underline{x}_n))\Big) = \kappa\big(d_{\mathrm{E}}(\underline{y}_m, \underline{y}_n)\big)$$

where $\underline{y}_m = f_i(\underline{x}_m)$ and $\underline{y}_n = f_i(\underline{x}_n)$ are elements of $\mathbb{R}^P$. Then, by assumption that $\kappa$ is a positive semi-definite covariance function over Euclidean space, $K$ is positive semi-definite. $\qquad\square$

We'll now define pseudometrics $d_i$ and associated isometries $f_i$ for both the bounded real and categorical cases.

## 2 Bounded Real Dimensions

Let's first define the 'difference' function $d_i$ on $\mathcal{X}$ and the isometry $f_i$ from $(\mathcal{X}, d_i)$ to $\mathbb{R}^2$, $d_{\mathrm{E}}$ for the case that the input $\mathcal{X}_i = [l_i, u_i]$ is bounded real; to emphasize that we're in this real case, we explicitly denote these functions as $d_i^{\mathrm{r}}$ and $f_i^{\mathrm{r}}$. We first define $d_i^{\mathrm{r}}$, recalling that $\delta_i(\underline{x})$ is true iff dimension $i$ is active given the instantiation of $i$'s ancestors in $\underline{x}$.

$$d_i^{\mathrm{r}}(\underline{x}, \underline{x}') = \begin{cases} 0 & \text{if } \delta_i(\underline{x}) = \delta_i(\underline{x}') = \text{false} \\ w_i & \text{if } \delta_i(\underline{x}) \neq \delta_i(\underline{x}') \\ w_i \sqrt{2} \sqrt{1 - \cos(\pi \frac{x_i - x_i'}{u_i - l_i})} & \text{if } \delta_i(\underline{x}) = \delta_i(\underline{x}') = \text{true.} \end{cases}$$

As desired, if $i$ is inactive for both $\underline{x}$ and $\underline{x}'$, $d_i^{\mathrm{r}}$ specifies that $g(\underline{x})$ and $g(\underline{x}')$ should not differ owing to differences between $x_i$ and $x_i'$.

> FH: Unfortunately, intuitively this isn't quite as nice as you wrote since the second case is not the one maximizing difference (it's up to $\sqrt{2}$ smaller than the third case).

> MO: Ah, you're right: I forgot that we'd left it like that. I think we could define the third case as $\frac{w_i}{\sqrt{2}}\sqrt{1 - \cos(\pi \frac{x_i - x_i'}{u_i - l_i})}$, and replace the $w_i$ in the isometry definition with $\frac{w_i}{2}$, if we wanted to and everything would still work. What would make more sense for the application, do you think?

If $i$ is inactive for exactly one of $\underline{x}$ and $\underline{x}'$, $g(\underline{x})$ and $g(\underline{x}')$ are almost as different as is possible with respect to dimension $i$, regardless of the actual values of $x_i$ and $x_i'$.[1] Finally, if $i$ is active for both $\underline{x}$ and $\underline{x}'$, the difference between $g(\underline{x})$ and $g(\underline{x}')$ due to $x_i$ and $x_i'$ increases monotonically with increasing $|x_i - x_i'|$.

---

[1] Note that $\underline{x}$ and $\underline{x}'$ must differ in at least one ancestor dimension of $i$ in order for $\delta_i(\underline{x}) \neq \delta_i(\underline{x}')$ to hold, such that differences in the activity of dimension $i$ are penalized both in distance $d_i$ and in the distance for the ancestor dimension causing the difference in $i$'s activity.

.

**Proposition 3.** $d_i^{\mathrm{r}}$ *is a pseudometric on* $\mathcal{X}$.

*Proof.* The non-negativity and symmetry of $d_i^{\mathrm{r}}$ are trivially proven. To prove the triangle inequality, consider $\underline{x}, \underline{x}', \underline{x}'' \in \mathcal{X}$.

Case 1: $\delta_i(\underline{x}) = \delta_i(\underline{x}') =$ false, such that $d_i^{\mathrm{r}}(\underline{x}, \underline{x}') = 0$. Here, from non-negativity, clearly $d_i^{\mathrm{r}}(\underline{x}, \underline{x}') = 0 \leq d_i^{\mathrm{r}}(\underline{x}, \underline{x}'') + d_i^{\mathrm{r}}(\underline{x}', \underline{x}'')$.

Case 2: $\delta_i(\underline{x}) \neq \delta_i(\underline{x}')$, such that such that $d_i^{\mathrm{r}}(\underline{x}, \underline{x}') = w_i$. Without loss of generality, assume $\delta_i(\underline{x}) =$ true, $\delta_i(\underline{x}') =$ false and $\delta_i(\underline{x}'') =$ true.

$$d_i^{\mathrm{r}}(\underline{x}, \underline{x}'') + d_i^{\mathrm{r}}(\underline{x}', \underline{x}'') = d_i^{\mathrm{r}}(\underline{x}, \underline{x}'') + w_i \tag{6}$$

Hence $d_i^{\mathrm{r}}(\underline{x}, \underline{x}'') + d_i^{\mathrm{r}}(\underline{x}', \underline{x}'') \geq w_i = d_i^{\mathrm{r}}(\underline{x}, \underline{x}')$ by non-negativity.

Case 3: $\delta_i(\underline{x}) = \delta_i(\underline{x}') =$ true, such that $d_i^{\mathrm{r}}(\underline{x}, \underline{x}') = w_i\sqrt{2}\sqrt{1 - \cos(\pi \frac{x_i - x_i'}{u_i - l_i})}$. If $\delta_i(\underline{x}'') =$ false,

$$d_i^{\mathrm{r}}(\underline{x}, \underline{x}'') + d_i^{\mathrm{r}}(\underline{x}', \underline{x}'') = 2w_i \geq w_i\sqrt{2}\sqrt{1 - \cos(\pi \frac{x_i - x_i'}{u_i - l_i})} = d_i^{\mathrm{r}}(\underline{x}, \underline{x}'). \tag{7}$$

If $\delta_i(\underline{x}'') = $ true, consider the worst possible case in which, without loss of generality, $x_i = l_i$ and $x_i' = u_i$, such that $d_i^{\mathrm{r}}(\underline{x}, \underline{x}') = 2w_i^2$. We define the abbreviation $\beta'' = \frac{x_i'' - l_i}{u_i - l_i}$, giving

$$
\begin{aligned}
\big(d_i^{\mathrm{r}}(\underline{x}, \underline{x}'') + d_i^{\mathrm{r}}(\underline{x}', \underline{x}'')\big)^2 &= 2w_i^2 \left( \sqrt{1 - \cos(\pi\beta'')} + \sqrt{1 - \cos\big(\pi(1 - \beta'')\big)} \right)^2 \\
&= 2w_i^2 \bigg( 2 - \cos(\pi\beta'') - \cos\big(\pi(1 - \beta'')\big) \\
&\qquad\quad + 2\sqrt{\Big(1 - \cos(\pi\beta'')\Big)\Big(1 - \cos\big(\pi(1 - \beta'')\big)\Big)} \bigg) \\
&= 2w_i^2 \bigg( 2 + 2\sqrt{1 + \cos(\pi\beta'')\cos\big(\pi(1 - \beta'')\big)} \bigg) \\
&= 4w_i^2 \big( 1 + |\sin \pi\beta''| \big) \\
&\geq 4w_i^2 = d_i^{\mathrm{r}}(\underline{x}, \underline{x}')^2.
\end{aligned}
\tag{8}
$$

Hence, from non-negativity, we have $d_i^{\mathrm{r}}(\underline{x}, \underline{x}'') + d_i^{\mathrm{r}}(\underline{x}', \underline{x}'') \geq d_i^{\mathrm{r}}(\underline{x}, \underline{x}')$. $\qquad\square$

Now we define an isometric embedding $f_i^{\mathrm{r}}$ of $(\mathcal{X}, d_i^{\mathrm{r}})$ into $(\mathbb{R}^2, d_E)$:

$$
f_i^{\mathrm{r}}(\underline{x}) = \begin{cases} [0, 0]^\mathsf{T} & \text{if } \delta_i(\underline{x}) = \text{ false} \\ w_i[\sin \pi \frac{x_i}{u_i - l_i}, \cos \pi \frac{x_i}{u_i - l_i}]^\mathsf{T} & \text{otherwise.} \end{cases}
$$

**Proposition 4.** $f_i^{\mathrm{r}}$ *is an isometry from* $(\mathcal{X}, d_i^{\mathrm{r}})$ *to* $(\mathbb{R}^2, d_E)$.

*Proof.* Consider two inputs $\underline{x}, \underline{x}' \in \mathcal{X}$. We need to show that $d_E\big(f_i^{\mathrm{r}}(\underline{x}), f_i^{\mathrm{r}}(\underline{x}')\big) = d_i^{\mathrm{r}}(\underline{x}, \underline{x}')$. We use the abbreviation $\alpha = \pi \frac{x_i}{u_i - l_i}$ and $\alpha' = \pi \frac{x_i'}{u_i - l_i}$ and consider the following three possible cases of dimension $i$ being active or inactive in $\underline{x}$ and $\underline{x}'$.

Case 1: $\delta_i(\underline{x}) = \delta_i(\underline{x}') = $ false. In this case, we trivially have

$$
d_E(f_i^{\mathrm{r}}(\underline{x}), f_i^{\mathrm{r}}(\underline{x}')) = d_E([0, 0]^\mathsf{T}, [0, 0]^\mathsf{T}) = 0 = d_i^{\mathrm{r}}(\underline{x}, \underline{x}').
$$

Case 2: $\delta_i(\underline{x}) \neq \delta_i(\underline{x}')$. In this case, we have

$$
d_E(f_i^{\mathrm{r}}(\underline{x}), f_i^{\mathrm{r}}(\underline{x}')) = d_E([\sin \alpha, \cos \alpha]^\mathsf{T}, [0, 0]^\mathsf{T}) = \sqrt{w_i^2(\sin^2 \alpha + \cos^2 \alpha)} = w_i = d_i^{\mathrm{r}}(\underline{x}, \underline{x}'),
$$

and symmetrically for $d_E([0, 0]^\mathsf{T}, [\sin \alpha, \cos \alpha]^\mathsf{T})$.

Case 3: $\delta_i(\underline{x}) = \delta_i(\underline{x}') =$ true. We have:

$$
\begin{aligned}
d_{\mathrm{E}}(f_i^{\,\mathrm{r}}(\underline{x}), f_i^{\,\mathrm{r}}(\underline{x}')) &= d_{\mathrm{E}}(w_i[\sin\alpha, \cos\alpha]^\top, w_i[\sin\alpha', \cos\alpha']^\top) \\
&= w_i\sqrt{(\sin\alpha - \sin\alpha')^2 + (\cos\alpha - \cos\alpha')^2} \\
&= w_i\sqrt{\sin^2\alpha - 2\sin\alpha\sin\alpha' + \sin^2\alpha' + \cos^2\alpha - 2\cos\alpha\cos\alpha' + \cos^2\alpha'} \\
&= w_i\sqrt{(\sin^2\alpha + \cos^2\alpha) + (\sin^2\alpha' + \cos^2\alpha') - 2(\sin\alpha\sin\alpha' + \cos\alpha\cos\alpha')} \\
&= w_i\sqrt{1 + 1 - 2\cos(\alpha - \alpha')} \\
&= w_i\sqrt{2}\sqrt{1 - \cos(\pi\frac{x_i - x_i'}{u_i - l_i})} = d_i^{\,\mathrm{r}}(\underline{x}, \underline{x}'),
\end{aligned}
\tag{9}
$$

where (9) follows from the previous line by using the identity

$$
\cos(a - b) = \cos a \cos b + \sin a \sin b.
$$

$\square$

# 3   Categorical Dimensions

Now let's define $f_i^{\mathrm{c}}$ and $d_i^{\mathrm{c}}$ for the case that the input $\mathcal{X}_i = \{v_{i,1}, \ldots, v_{i,m_i}\}$ is categorical with $m_i$ possible values. Proceeding as above, we define a pseudometric $d_i^{\mathrm{c}}$ on $\mathcal{X}$ and an isometry from $(\mathcal{X}, d_i^{\mathrm{c}})$ to $(\mathbb{R}^{m_i}, d_{\mathrm{E}}^{m_i})$.

Firstly,

$$
d_i^{\mathrm{c}}(\underline{x}, \underline{x}') = \begin{cases} 0 & \text{if } \delta_i(\underline{x}) = \delta_i(\underline{x}') = \text{false} \\ w_i & \text{if } \delta_i(\underline{x}) \neq \delta_i(\underline{x}') \\ w_i\sqrt{2}\,\mathbb{I}_{x_i \neq x_i'} & \text{if } \delta_i(\underline{x}) = \delta_i(\underline{x}') = \text{true.} \end{cases}
$$

**Proposition 5.** $d_i^{\mathrm{c}}$ *is a pseudometric on* $\mathcal{X}$.

*Proof.* A trivial modification of the proof to Proposition 3. $\square$

Secondly,

$$
f_i^{\mathrm{c}}(\underline{x}) = \begin{cases} \underline{0} \in \mathbb{R}^{m_i} & \text{if } \delta_i(\underline{x}) = \text{false} \\ w_i\,\underline{e_j} & \delta_i(\underline{x}) = \text{true and } x_i = v_{i,j}, \end{cases}
$$

where $\underline{e_j} \in \mathbb{R}^{m_i}$ is zero in all dimensions except $j$, where it it 1.

**Proposition 6.** $f_i^{\mathrm{c}}$ *is an isometry from* $(\mathcal{X}, d_i^{\mathrm{c}})$ *to* $(\mathbb{R}^{m_i}, d_E^{m_i})$.

*Proof.* Consider two inputs $\underline{x}, \underline{x}' \in \mathcal{X}$. As in the proof of Proposition 4, we need to show that $d_i^{\mathrm{c}}(\underline{x}, \underline{x}') = d_{\mathrm{E}}^{m_i}(f_i^{\mathrm{c}}(\underline{x}), f_i^{\mathrm{c}}(\underline{x}'))$ and consider the following cases.

Case 1: $\delta_i(\underline{x}) = \delta_i(\underline{x}') =$ false. In this case, we trivially have

$$
d_{\mathrm{E}}^{m_i}(f_i^{\,\mathrm{r}}(\underline{x}), f_i^{\,\mathrm{r}}(\underline{x}')) = d_{\mathrm{E}}^{m_i}(\underline{0}, \underline{0}) = 0 = d_i^{\,\mathrm{r}}(\underline{x}, \underline{x}').
$$

Case 2: $\delta_i(\underline{x}) \neq \delta_i(\underline{x}')$. In this case, we have

$$d_{\mathrm{E}}^{m_i}(f_i^{\mathrm{c}}(\underline{x}), f_i^{\mathrm{c}}(\underline{x}')) = d_{\mathrm{E}}^{m_i}(w_i\,\underline{e_j}, \underline{0}) = w_i = d_i(\underline{x}, \underline{x}'),$$

and symmetrically for $d_{\mathrm{E}}(\underline{0}, w_i\,\underline{e_j})$.

Case 3: $\delta_i(\underline{x}) = \delta_i(\underline{x}') = \mathrm{true}$. If $x_i = x_i' = v_{i,j}$, we have

$$d_{\mathrm{E}}^{m_i}(f_i^{\mathrm{c}}(\underline{x}), f_i^{\mathrm{c}}(\underline{x}')) \;\;=\;\; d_{\mathrm{E}}^{m_i}(w_i\,\underline{e_j}, w_i\,\underline{e_j}) = 0 = d_i^{\mathrm{c}}(\underline{x}, \underline{x}')$$

If $x_i = v_{i,j} \neq v_{i,j'} = x_i' =$, we have

$$d_{\mathrm{E}}(f_i^{\mathrm{c}}(\underline{x}), f_i^{\mathrm{c}}(\underline{x}')) \;\;=\;\; d_{\mathrm{E}}^{m_i}(w_i\,\underline{e_j}, w_i\underline{e_{j'}}) = w_i\sqrt{2} = d_i^{\mathrm{c}}(\underline{x}, \underline{x}')$$

$\square$