

An equivalent formulation of matrix kernels

Miaoyan Wang, July 30, 2020

For ease of notation, we assume the feature matrices are symmetric in that $\mathbf{X} = \mathbf{X}^T$ and $d_1 = d_2 = d$. The adaptation to non-symmetric matrices are easy to derive.

- Let $\text{Sym}_d(\mathbb{R}) = \{\mathbf{X} \mid \mathbf{X} = \llbracket x_{ij} \rrbracket, x_{ij} = x_{ji} \in \mathbb{R}, \text{ for all } (i, j) \in [d] \times [d]\}$ denote the collection of d -by- d symmetric matrices with each entry taking values in \mathbb{R} .
- Let \mathcal{H} denote a possibly infinite dimensional Hilbert space.
- Let $\mathcal{H}^{d \times d} = \{\mathbf{X} \mid \mathbf{X} = \llbracket x_{ij} \rrbracket, x_{ij} \in \mathcal{H}, \text{ for all } (i, j) \in [d] \times [d]\}$ denote the collection of d -by- d matrices with each entry taking value in \mathcal{H} . Similarly, $\text{Sym}_d(\mathcal{H})$ denotes the collection of d -by- d symmetric matrices defined over \mathcal{H} .

Matrix algebraic operations are carried over from $\mathbb{R}^{d \times d}$ to $\mathcal{H}^{d \times d}$. The main difference is that the multiplication in \mathbb{R} is replaced by inner product in \mathcal{H} .

Proposition 1. Let $\mathbf{B} = \llbracket b_{ij} \rrbracket$ and $\mathbf{B}' = \llbracket b'_{ij} \rrbracket$ be two matrices in $\mathcal{H}^{d \times d}$. Let $\mathbf{P} = \llbracket p_{ij} \rrbracket \in \mathbb{R}^{d \times r}$ be a real-valued matrix.

1. Sum: $\mathbf{B} + \mathbf{B}' = \llbracket b_{ij} + b'_{ij} \rrbracket \in \mathcal{H}^{d \times d}$.
2. Inner product: $\langle \mathbf{B}, \mathbf{B}' \rangle = \sum_{ij} \langle b_{ij}, b'_{ij} \rangle \in \mathbb{R}$.
3. Linear combination: $\mathbf{B}\mathbf{P} = \llbracket c_{ij} \rrbracket \in \mathcal{H}^{d \times r}$, where $c_{ij} = \sum_{s \in [d]} b_{is} p_{sj} \in \mathcal{H}^2$ for all $(i, j) \in [d] \times [r]$.
4. Matrix product: $\mathbf{B}\mathbf{B}' = \llbracket c_{ij} \rrbracket \in \mathcal{H}^{d \times d}$, where $c_{ij} = \sum_{s \in [d]} \langle b_{is}, b'_{sj} \rangle$.

Now we are ready to present the matrix kernel and associated feature mapping.

First viewpoint: feature mapping

Let $\phi: \mathbb{R}^d \rightarrow \mathcal{H}$ be a feature mapping associated with a classical kernel $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$. Define a feature mapping Φ over the matrix space

$$\Phi: \text{Sym}_d(\mathbb{R}) \rightarrow \text{Sym}_d(\mathcal{H}^2)$$

$$\mathbf{X} \mapsto \Phi(\mathbf{X}) = \llbracket \Phi(\mathbf{X})_{ij} \rrbracket, \quad \begin{array}{l} \text{special case:} \\ [\Phi(\mathbf{X})]_{ij} = (\phi(\mathbf{X}_{i,:}), \phi(\mathbf{X}_{j,:})) \rightarrow \text{in your current draft} \\ \text{asymmetric case:} \\ \text{new feature} = \begin{bmatrix} \mathbf{X} & \mathbf{0} \\ \mathbf{0} & \mathbf{X} \end{bmatrix} \end{array}$$

where each element of $\Phi(\mathbf{X})$ is a pair of possibly infinite dimensional features

$$[\Phi(\mathbf{X})]_{ij} \stackrel{\text{def}}{=} \begin{cases} (\phi(\mathbf{X}_{i,:}), \phi(\mathbf{X}_{j,:})), & \text{if } i \geq j, \\ [\Phi(\mathbf{X})]_{ji}, & \text{if } i < j. \end{cases}$$

Note that the entry f_{ij} takes value in \mathcal{H}^2 for all $(i, j) \in [d] \times [d]$. Furthermore, $\Phi(\mathbf{X})$ is a symmetric matrix in the sense that $[\Phi(\mathbf{X})]_{ij} = [\Phi(\mathbf{X})]_{ji}$ for all $(i, j) \in [d] \times [d]$. We propose to consider decision



functions $f: \text{Sym}_d(\mathbb{R}) \rightarrow \mathbb{R}$ using the linear functions with respect to $\Phi(\mathbf{X}) \in \text{Sym}_d(\mathcal{H}^2)$,

$$f(\mathbf{X}) \stackrel{\text{def}}{=} \langle \mathbf{B}, \Phi(\mathbf{X}) \rangle, \text{ where } \mathbf{B} \in \text{Sym}_d(\mathcal{H}^2) \text{ and } \text{rank}(\mathbf{B}) \leq r. \quad (1)$$

Here the parameter matrix $\mathbf{B} = \llbracket b_{ij} \rrbracket$ is a d -by- d symmetric matrix with entries defined in \mathcal{H}^2 . Suppose \mathbf{B} admits low-rank decomposition, $\mathbf{B} = \mathbf{P}^T \mathbf{C} \mathbf{P}$. The class of functions (1) induced by all possible low-rank \mathbf{B} is

$$\begin{aligned} \mathcal{F} &= \left\{ f: \mathbf{X} \mapsto \langle \mathbf{C}, \mathbf{P}^T \Phi(\mathbf{X}) \mathbf{P} \rangle \mid \mathbf{P} \mathbf{P}^T = \mathbf{I}, \mathbf{P} \in \mathbb{R}^{d \times r}, \mathbf{C} \in \text{Sym}_d(\mathcal{H}^2) \right\} \\ &= \left\{ f \in \text{RKHS generated by } \mathcal{K}(\mathbf{P}) \mid \mathbf{P} \mathbf{P}^T = \mathbf{I}, \mathbf{P} \in \mathbb{R}^{d \times r} \right\} \\ &= \left\{ f \in \text{RKHS generated by } \mathcal{K}(\mathbf{W}) \mid \mathbf{W} \succeq 0, \text{rank}(\mathbf{W}) \leq r \right\}, \end{aligned}$$

where $\mathbf{W} = \mathbf{P}^T \Lambda^2 \mathbf{P}$ for some positive definite diagonal matrix $\Lambda \in \mathbb{R}^{r \times r}$ (see below).

Second viewpoint: matrix kernel

Definition 1 (Kernel defined in matrix space). Let $K(\cdot, \cdot)$ be a classical kernel defined in vector space \mathbb{R}^d , and $\mathbf{W} = \llbracket w_{ij} \rrbracket \in \mathbb{R}^{d \times d}$ be a rank- r semi-positive definite matrix. Then \mathbf{K} and \mathbf{W} induce a matrix kernel \mathcal{K} :

call this K as K_col

$$\begin{aligned} \mathcal{K}: \text{Sym}_d(\mathbb{R}) \times \text{Sym}_d(\mathbb{R}) &\mapsto \mathbb{R}, \\ (\mathbf{X}, \mathbf{X}') &\mapsto \mathcal{K}(\mathbf{X}, \mathbf{X}') = \sum_{i,j \in [d]} w_{ij} K(\mathbf{x}_i, \mathbf{x}'_j), \end{aligned}$$

conjecture: r_col K_col(X,X') + r_row K_row(X,X')

where $\mathbf{x}_i, \mathbf{x}'_j$ denote the i -th and j -th columns of \mathbf{X}, \mathbf{X}' , respectively.

Oftentime, the generator kernel K is specified by users, whereas the projection kernel \mathbf{W} is learned by our algorithm. In particular $\mathbf{W} \propto \mathbf{I}_{d \times d}$ corresponds to the classical SVM. We use $\mathcal{K} = \mathcal{K}(\mathbf{W})$ to denote the reproducing kernel hilbert space (RKHS) induced by \mathbf{W} .

Connection to our learning algorithm

We consider the optimization over the union of RKHS induced by low rank \mathbf{W} :

$$\max_{f \in \mathcal{F}(r)} L(f) = \max_{\substack{\text{rank}(\mathbf{W}) \leq r, \\ \mathbf{W} \succeq 0}} \max_{f \in \text{RKHS}(\mathcal{K}(\mathbf{W}))} L(f).$$