

Kernel machines

Nathanaël Carraz Rakotonirina

Mathématiques Informatique et Statistique Appliquées (MISA)
Université d'Antananarivo

Motivations

- ▶ We want to use features that are more appropriate (instead of just the raw inputs) for a given problem.
- ▶ Instead of operating on the inputs \mathbf{x} , we operate on features $\phi(\mathbf{x})$ (using the feature mapping ϕ) which can result in non-linear models with more capacity.
- ▶ We want to use these features efficiently in our models.

Kernel trick

Given a feature mapping ϕ , we define the corresponding **kernel** to be

$$\begin{aligned}\mathcal{K} : \mathcal{X} \times \mathcal{X} &\rightarrow \mathbb{R}^+ \\ (\mathbf{x}, \mathbf{x}') &\mapsto \mathcal{K}(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')\end{aligned}$$

The **kernel trick** is about replacing the dot products $\mathbf{x}^\top \mathbf{x}'$ in our model with the kernel function $\mathcal{K}(\mathbf{x}, \mathbf{x}')$. By doing so

- ▶ The model would now be learning using the features ϕ .
- ▶ There is a way to efficiently calculate $\mathcal{K}(\mathbf{x}, \mathbf{x}')$ without having to explicitly find and compute the feature vectors $\phi(\mathbf{x})$ (which can be very expensive).

How to find a kernel that is valid or corresponds to some feature mapping ϕ ?

Mercer kernel

A **Mercer kernel** or **positive definite kernel** is a function $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^+$ such that:

- ▶ It is symmetric $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \mathcal{K}(\mathbf{x}', \mathbf{x})$
- ▶ For any set of (unique) points $\{\mathbf{x}_i\}_{i=1}^N$, and any numbers $c_i \in \mathbb{R}$

$$\sum_{i=1}^N \sum_{j=1}^N \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) c_i c_j \geq 0$$

There is another way to define it. Given a set of N points, the **Gram matrix** \mathbf{K} is an $N \times N$ symmetric matrix with entries $\mathbf{K}_{i,j} = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$. \mathcal{K} is a Mercer kernel iff the Gram matrix \mathbf{K} is positive definite for any set of (distinct) points $\{\mathbf{x}_i\}_{i=1}^N$.

Mercer theorem

A kernel $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^+$ can be computed by an inner product of some feature vectors iff for any set of points $\{\mathbf{x}_i\}_{i=1}^N$, it is positive definite.

Example of Mercer kernels

- ▶ Linear : $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$
- ▶ Quadratic : $\mathcal{K}(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + b)^2$
- ▶ Polynomial : $\mathcal{K}(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + b)^p$
- ▶ Gaussian (RBF) $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \exp \left(- \frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2} \right)$
- ▶ Laplacian : $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \exp \left(- \frac{\|\mathbf{x} - \mathbf{x}'\|}{\sigma} \right)$

Some kernels may contain hyperparameters that needs to be tuned using cross-validation.

Making new kernels

Given valid kernels $\mathcal{K}_1(\mathbf{x}, \mathbf{x}')$ and $\mathcal{K}_2(\mathbf{x}, \mathbf{x}')$, we can create a new kernel using:

- ▶ $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \mathcal{K}_1(\mathbf{x}, \mathbf{x}') + \mathcal{K}_2(\mathbf{x}, \mathbf{x}')$
- ▶ $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \mathcal{K}_1(\mathbf{x}, \mathbf{x}') \times \mathcal{K}_2(\mathbf{x}, \mathbf{x}')$
- ▶ $\mathcal{K}(\mathbf{x}, \mathbf{x}') = c\mathcal{K}_1(\mathbf{x}, \mathbf{x}')$ for any constant $c > 0$
- ▶ $\mathcal{K}(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})\mathcal{K}_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$ for any function f
- ▶ $\mathcal{K}(\mathbf{x}, \mathbf{x}') = q(\mathcal{K}_1(\mathbf{x}, \mathbf{x}'))$ for any function polynomial q with non-negative coefficients
- ▶ $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \exp(\mathcal{K}_1(\mathbf{x}, \mathbf{x}'))$
- ▶ $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{A} \mathbf{x}'$ for any positive semi-definite matrix \mathbf{A}

Explore further

- ▶ Other examples of kernels
- ▶ Kernels for structured inputs (strings, time series, graphs, images)
- ▶ Kernel PCA
- ▶ Kernel ridge regression
- ▶ Automatically Choosing a Kernel