# Statistical Machine Learning
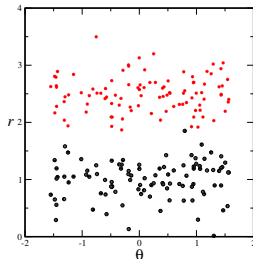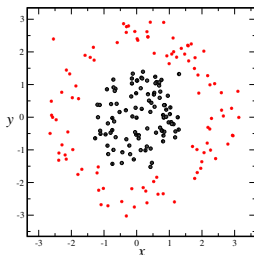
## Lecture 9: Feature Space and Kernels

2022-23

# Feature Space

▶ Some time changing data representation $x \in \mathcal{X} \subseteq \mathbb{R}^d$ is a fact of convenience, $\phi : \mathcal{X} \to \mathcal{F} \subseteq \mathbb{R}^M$.

▶ If $M < d$ there is a dimensionality reduction (like in PCA).

▶ It may simplify the classification boundary.

# Mapping into Feature Space

- In the previous slide the transformation is
  $\phi : \mathscr{X} \subseteq \mathbb{R}^2 \to \mathscr{F} \subseteq (-\pi, \pi] \times \mathbb{R}^+$, with
  $\theta = \phi_1(\boldsymbol{x}) = \arctan(y/x)$ and $r = \phi_2(\boldsymbol{x}) = \|\boldsymbol{x}\|_2$ .
- A non-linear mapping $(\phi)$ can be used in order to apply a linear machine in the feature space:, i.e. given
  $\phi : \mathscr{X} \subseteq \mathbb{R}^d \to \mathscr{F} \subseteq \mathbb{R}^M$, $\boldsymbol{w} \in \Omega \subseteq \mathbb{R}^M$ and $w_0 \in \mathbb{R}$ :

$$f(\boldsymbol{x}) = \boldsymbol{w}^T \phi(\boldsymbol{x}) + w_0.$$

- We can also represent this linear machine in the dual representation:

$$f(\boldsymbol{x}) = \sum_{n=1}^{N} \alpha_n y_n \phi(\boldsymbol{x}_n)^T \phi(\boldsymbol{x}) + w_0.$$

.

# Mapping into Feature Space

- It is frequently used the bra-ket notation:

$$\phi(x) = |\phi(x)\rangle \ \text{column vector,}$$
$$\phi(x)^T = \langle\phi(x)| \ \text{row vector,}$$
$$\phi(x)^T\phi(y) = \langle\phi(x), \phi(y)\rangle \ \text{inner product.}$$

- The dual representation becomes:

$$f(x) = \sum_{n=1}^{N} \alpha_n y_n \langle\phi(x_n), \phi(x)\rangle + w_0.$$

- If we have a means to compute the inner product $\langle\phi(x_n), \phi(x)\rangle$ in feature space directly as a function of the original input points, it becomes possible to merge the two steps needed to build a non-linear learning machine.

.

# Kernels

- A *kernel* is a function $K : \mathscr{X}^2 \to \mathbb{R}$, $\mathscr{X}$ an inner product space, that for all $x, z \in \mathscr{X}$ :

$$K(x, z) = \langle \phi(x), \phi(z) \rangle ,$$

  where $\phi : \mathscr{X} \to \mathscr{F}$.

- The use of kernels make the representation into the feature space implicit. The feature vectors $\phi(x)$ do not need to be explicitly computed.

- The dual representation induces that the training examples appear in the form of inner products, stored into a Gramm matrix $G$ (symmetric and positive definite).

# Kernels

- Equivalently, the inner product between feature vectors involving the elements of the data set, can be stored in a kernel matrix:
$$[K]_{m,n} = K(x_m, x_n).$$

- The kernel function is also used to represent the linear machine:
$$f(x) = \sum_{n=1}^{N} \alpha_n y_n K(x_n, x) + w_0.$$

- Observe that by the use of an appropriate kernel, knowledge about the feature map $\phi$ is not needed.

.

# Examples of Kernels

- The identity:
$$K(\boldsymbol{x}, \boldsymbol{y}) = \langle \boldsymbol{x}, \boldsymbol{y} \rangle.$$

- Linear: given $\boldsymbol{A} \in \mathbb{R}^{m \times d}$

$$K(\boldsymbol{x}, \boldsymbol{y}) = \langle \boldsymbol{A}\boldsymbol{x}, \boldsymbol{A}\boldsymbol{y} \rangle = (\boldsymbol{A}\boldsymbol{x})^T \boldsymbol{A}\boldsymbol{y} = \boldsymbol{x}^T \boldsymbol{A}^T \boldsymbol{A}\boldsymbol{y} = \boldsymbol{x}^T \boldsymbol{B}\boldsymbol{y},$$

with $\boldsymbol{B} = \boldsymbol{A}^T \boldsymbol{A}$.

# Examples of Kernels

- Square:

$$K(\boldsymbol{x}, \boldsymbol{y}) = \langle \boldsymbol{x}, \boldsymbol{y} \rangle^2 = \sum_{i=1}^{d} \sum_{j=1}^{d} x_i x_j y_i y_j = \sum_{(i,j)=(1,1)}^{(d,d)} (x_i x_j)(y_i y_j),$$

  which is equivalent to the inner product between the feature vectors:

$$[\phi(\boldsymbol{x})]_{(i,j)} = (x_i x_j).$$

- Generalized square: $(\langle \boldsymbol{x}, \boldsymbol{y} \rangle + c)^2$. And other powers.

# Properties of Kernels

- Symmetry:

$$K(\boldsymbol{x}, \boldsymbol{y}) = \langle \phi(\boldsymbol{x}), \phi(\boldsymbol{y}) \rangle = \langle \phi(\boldsymbol{y}), \phi(\boldsymbol{x}) \rangle = K(\boldsymbol{y}, \boldsymbol{x}).$$

- Cauchy-Schwarz inequality:

$$\begin{aligned}
(K(\boldsymbol{x}, \boldsymbol{y}))^2 &= \langle \phi(\boldsymbol{x}), \phi(\boldsymbol{y}) \rangle^2 \\
&\leq \|\phi(\boldsymbol{x})\|_2^2 \|\phi(\boldsymbol{y})\|_2^2 = \langle \phi(\boldsymbol{x}), \phi(\boldsymbol{x}) \rangle \langle \phi(\boldsymbol{y}), \phi(\boldsymbol{y}) \rangle \\
&\leq K(\boldsymbol{x}, \boldsymbol{x}) K(\boldsymbol{y}, \boldsymbol{y})
\end{aligned}$$

# Properties of Kernels

- Suppose that the input space $\mathscr{X} = \{x_1, \ldots, x_M\}$ is finite ($M < \infty$). Suppose $K : \mathscr{X} \times \mathscr{X} \to \mathbb{R}$ is symmetric. Consider the matrix $K$ with entries $[K]_{i,j} = K(x_i, x_j)$.

- Because of the symmetry there exist $\Lambda, V \in \mathbb{R}^{M \times M}$ such that $VV^T = V^TV = I$ and $\Lambda$ diagonal, and $\mathrm{diag}(\Lambda) = \{\lambda_1, \ldots, \lambda_M\}$ eigenvalues of $K$ with eigenvectors given by the columns of $V$.

- Assume that $\lambda_1 \geq \cdots \geq \lambda_M \geq 0$. Let us define the matrix $X \in \mathbb{R}^{M \times M}$ with columns $x_k$. Consider the transformation $\phi(X) = \Lambda^{1/2}V^T$ such that $\phi(X)^T\phi(X) = K$.

# Properties of Kernels

- In particular $\phi(x_k) = \operatorname{col}(\boldsymbol{\Lambda}^{1/2}\boldsymbol{V}^T)_k$ and $\phi(x_j)^T = \operatorname{row}(\boldsymbol{V}\boldsymbol{\Lambda}^{1/2})_j$, therefore $\langle\phi(x_j), \phi(x_k)\rangle = \sum_{\ell=1}^{M} \lambda_\ell v_{\ell,j} v_{\ell,k} = K(x_j, x_k)$.

- Let as assume now that the feature space $\mathscr{F}$ is infinite Hilbert space (a space that is complete and its norm is given by an inner product) (remember $\phi : \mathscr{X} \to \mathscr{F}$) we can generalize the definition of the inner product as:

$$\langle\phi(x), \phi(z)\rangle = \sum_{\ell=1}^{\infty} \lambda_\ell \phi_\ell(x)\phi_\ell(z) = K(x, z),$$

with $\lambda_1 \geq \lambda_2 \geq \cdots \geq 0$. (Mercer's Theorem).

# Properties of Kernels

- Observe that this definition requires that for all $\psi \in \mathscr{F}$ $\sum_{\ell=1}^{\infty} \lambda_\ell \psi_\ell^2 < \infty$, so $\mathscr{F}$ is the space that contains all sequences $\psi = (\psi_1, \psi_2, \dots)$ with finite inner product norm given by $\{\lambda_\ell\}_{\ell=1}^{\infty}$.
- The linear machine will be represented by:

$$f(\mathbf{x}) = \sum_{\ell=1}^{\infty} \lambda_\ell \psi_\ell \phi_\ell(\mathbf{x}) + w_0 = \sum_{n=1}^{N} \alpha_n t_n K(\mathbf{x}_n, \mathbf{x}) + w_0,$$

where

$$\psi = \sum_{n=1}^{N} \alpha_n t_n \phi(\mathbf{x}_n).$$

# Mercer's Theorem

- Theorem: Let $\mathscr{X}$ be a compact subset of $\mathbb{R}^d$ (compact is similar to require that all sequences of elements in $\mathscr{X}$ converge in $\mathscr{X}$.) Suppose $K$ is a continuous symmetric function such that the integral operator $T_K : L_2(\mathscr{X}) \to L_2(\mathscr{X})$,

$$(T_K f)(\boldsymbol{y}) = \int_{\mathscr{X}} \mathrm{d}\boldsymbol{x} K(\boldsymbol{y}, \boldsymbol{x}) f(\boldsymbol{x}),$$

- is positive, i.e.

$$\int_{\mathscr{X} \times \mathscr{X}} \mathrm{d}\boldsymbol{x}\mathrm{d}\boldsymbol{z} K(\boldsymbol{x}, \boldsymbol{z}) f(\boldsymbol{x}) f(\boldsymbol{z}) \geq 0, \qquad \forall f \in L_2(\mathscr{X}).$$

.

# Mercer's Theorem

▶ Then we can expand $K(x, z)$ in a uniformly convergent series (on $\mathcal{X} \times \mathcal{X}$) in terms of $T_K$ eigenfunctions $\phi_j \in L_2(\mathcal{X})$, normalised in such a way that $\|\phi_j\|_{L_2} = 1$, and positive associated eigenvalues $\lambda_j > 0$,

$$K(x, z) = \sum_{j=1}^{\infty} \lambda_j \phi_j(x) \phi_j(z).$$

▶ Observe that we do not need to compute the Mercer's features $\phi(x) \in \mathcal{F}$. We only need to compute their inner products through $K$, remember that

$$f(x) = \sum_{\ell \in \text{SV}} \alpha_\ell t_\ell K(x_\ell, x).$$

# Reproducing Kernel Hilbert Spaces

- Let us define the set $\mathscr{H} = \{\sum_{j=1}^{\ell} \alpha_j K(x_j, \cdot) : \ell \in \mathbb{N}, x_j \in \mathscr{X}, \alpha_j \in \mathbb{R}, j = 1, \ldots, \ell\}$. (Observe these are functions $\mathscr{H} \ni h : \mathscr{X} \to \mathbb{R}$).

- Observe that if $f(\cdot), g(\cdot) \in \mathscr{H}$ and $a, b \in \mathbb{R}$ then $af(\cdot) + bg(\cdot) \in \mathscr{H}$. Thus $\mathscr{H}$ is a vector space.

- Suppose $f(\cdot) = \sum_{j=1}^{\ell_f} \alpha_j K(x_j, \cdot)$ and $g(\cdot) = \sum_{j=1}^{\ell_g} \beta_j K(z_j, \cdot)$ both in $\mathscr{H}$. We define the inner product of $\mathscr{H}$ as

$$\langle f, g \rangle_{\mathscr{H}} = \sum_{j=1}^{\ell_f} \sum_{k=1}^{\ell_g} \alpha_j \beta_k K(x_j, z_k)$$

$$= \sum_{j=1}^{\ell_f} \alpha_j g(x_j) = \sum_{k=1}^{\ell_g} \beta_k f(z_k).$$

# Reproducing Kernel Hilbert Spaces

▶ The kernel matrices are positive semidefinite, i.e. for any collection of input space vectors $\{x_j \in \mathscr{X}\}_{j=1}^N$, $K \in \mathbb{R}^{N \times N}$ defined as $[K]_{ij} = K(x_i, x_j)$ is symmetric, real, and therefore positive semidefinite. Thus

$$\langle f, f \rangle_{\mathscr{H}} = \sum_{j=1}^{\ell_f} \sum_{k=1}^{\ell_f} \alpha_j K(x_j, x_k) \alpha_k = \boldsymbol{\alpha}^T K \boldsymbol{\alpha} \geq 0.$$

▶ Observe also that $K(x, \cdot) \in \mathscr{H}$ (with $\ell_K = 1$, $\alpha_1 = 1$, and $x_1 = x$). Thus

$$\langle f, K(x, \cdot) \rangle_{\mathscr{H}} = \sum_{j=1}^{\ell_f} \alpha_j K(x_j, x) = f(x),$$

which is known as the *reproducing property* of the kernel.

# Reproducing Kernel Hilbert Spaces

- $\mathscr{H}$ is known as the *Reproducing Kernel Hilbert Space*.
- Learning implies to find a suitable function $f \in \mathscr{H}$ that separates classes. In the dual representation $f(\mathbf{x}) = \sum_{j \in \text{SV}} \alpha_j K(\mathbf{x}_j, \mathbf{x}) \in \mathscr{H}$.