# Statistical Learning and Machine Learning
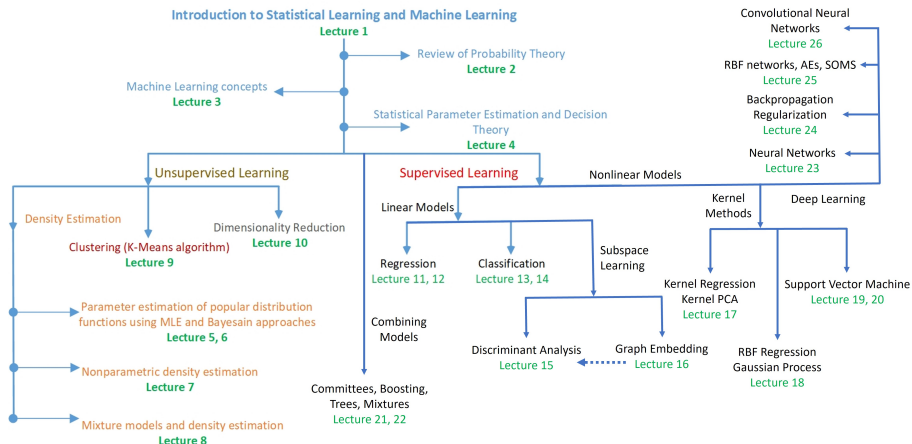## Lecture 17 - Kernel Methods 1

November 1, 2021

# Course overview and where do we stand



Introduction to Statistical Learning and Machine Learning
Lecture 1

Review of Probability Theory
Lecture 2

Machine Learning concepts
Lecture 3

Statistical Parameter Estimation and Decision Theory
Lecture 4

Convolutional Neural Networks
Lecture 26

RBF networks, AEs, SOMS
Lecture 25

Backpropagation
Regularization
Lecture 24

Neural Networks
Lecture 23

Unsupervised Learning

Supervised Learning

Nonlinear Models

Density Estimation

Clustering (K-Means algorithm)
Lecture 9

Dimensionality Reduction
Lecture 10

Linear Models

Kernel Methods

Deep Learning

Regression
Lecture 11, 12

Classification
Lecture 13, 14

Subspace Learning

Kernel Regression
Kernel PCA
Lecture 17

Support Vector Machine
Lecture 19, 20

Parameter estimation of popular distribution functions using MLE and Bayesian approaches
Lecture 5, 6

Combining Models

Nonparametric density estimation
Lecture 7

Discriminant Analysis
Lecture 15

Graph Embedding
Lecture 16

RBF Regression
Gaussian Process
Lecture 18

Mixture models and density estimation
Lecture 8

Committees, Boosting, Trees, Mixtures
Lecture 21, 22

# Why kernel methods?

Main idea in kernel methods:

- The prediction of the model is based on a linear combination of dot-products between data points:

$$\kappa(\mathsf{x}, \mathsf{x}') = \phi(\mathsf{x})^T \phi(\mathsf{x}') \tag{1}$$

where $\phi(\mathsf{x})$ is an (a-priori defined) (non)linear transformation of the data points $\mathsf{x}$.

- Observation: $\kappa(\cdot, \cdot)$ is a (scalar) value, and this is not depending on the dimensionality of $\mathsf{x}$. $\kappa(\cdot, \cdot)$ is called *kernel function*.

- Whenever a method can be expressed using (only) dot-products between data points, a so-called *dual representation* of the method is defined using $\kappa(\cdot, \cdot)$

- The dot-product between the data points can be replaced by any nonlinear function (satisfying the properties of kernel functions).

- Using a nonlinear kernel function, the nonlinear counterpart of the linear method is defined.

# Kernel Regression

Regularized sum-of-squares error:

$$\mathcal{J}(\mathsf{w}) = \frac{1}{2}\sum_{n=1}^{N}\left(\mathsf{w}^{T}\phi(\mathsf{x}_n) - t_n\right)^2 + \frac{\lambda}{2}\mathsf{w}^{T}\mathsf{w} \tag{2}$$

where $t_n$ is the target value for $\phi(\mathsf{x}_n)$ and $\lambda \geq 0$.

Using $\Phi = [\phi(\mathsf{x}_1)^{T}, \ldots, \phi(\mathsf{x}_N)^{T}] \in \mathbb{R}^{N \times |\phi|}$ and $\mathsf{t} = [t_1, \ldots, t_N]^{T} \in \mathbb{R}^{N}$:

$$\mathcal{J}(\mathsf{w}) = \frac{1}{2}\|\Phi\mathsf{w} - \mathsf{t}\|^2 + \frac{\lambda}{2}\mathsf{w}^{T}\mathsf{w} \tag{3}$$

# Kernel Regression

Setting $\theta \mathcal{J}(w)/\theta w = 0$, we get:

$$w = -\frac{1}{\lambda} \sum_{n=1}^{N} \left( w^T \phi(x_n) - t_n \right) \phi(x_n) = \sum_{n=1}^{N} \alpha_n \phi(x_n) = \Phi^T \boldsymbol{\alpha} \qquad (4)$$

where $\boldsymbol{\alpha} = [\alpha_1, \ldots, \alpha_N]^T$:

$$\alpha_n = -\frac{1}{\lambda} \left( w^T \phi(x_n) - t_n \right). \qquad (5)$$

Thus, we can obtain a *dual representation* for linear regression, leading to a non-linear version of regression.

# Kernel Regression

If we substitute $w = \Phi^T \boldsymbol{\alpha}$ in $\mathcal{J}(w)$:

$$\mathcal{J}(\boldsymbol{\alpha}) = \frac{1}{2}\boldsymbol{\alpha}^T \Phi \Phi^T \Phi \Phi^T \boldsymbol{\alpha} - \boldsymbol{\alpha}\Phi\Phi^T t + \frac{1}{2}t^T t + \frac{\lambda}{2}\boldsymbol{\alpha}^T \Phi \Phi^T \boldsymbol{\alpha} \quad (6)$$

We define the *Gram* matrix (or *kernel matrix*) $K = \Phi\Phi^T \in \mathbb{R}^{N \times N}$ having elements:

$$K_{nm} = \phi(x_n)^T \phi(x_m) = \kappa(x_n, x_m). \quad (7)$$

Using K, $\mathcal{J}(\boldsymbol{\alpha})$ becomes:

$$\mathcal{J}(\boldsymbol{\alpha}) = \frac{1}{2}\boldsymbol{\alpha}^T KK\boldsymbol{\alpha} - \boldsymbol{\alpha}^T Kt + \frac{1}{2}t^T t + \frac{\lambda}{2}\boldsymbol{\alpha}^T K\boldsymbol{\alpha} \quad (8)$$

Setting $\theta\mathcal{J}(\boldsymbol{\alpha})/\theta\boldsymbol{\alpha} = 0$:

$$\boldsymbol{\alpha} = (K + \lambda I_N)^{-1} t. \quad (9)$$

# Kernel Regression

For a data point $x_*$ we have:

$$y(x_*) = w^T \phi(x_*) = \boldsymbol{\alpha}^T \Phi \phi(x_*) = k(x_*)^T (K + \lambda I_N)^{-1} t \qquad (10)$$

where $k(x_*) \in \mathbb{R}^N$ has elements equal to $k_n(x_*) = \kappa(x_n, x_*)$.

Note that above we used $|\phi|$ to denote the dimensionality of the vectors $\phi(x_n)$. Since $|\phi|$ is 'hidden' in the dot-product, it can take any value (even infinity).
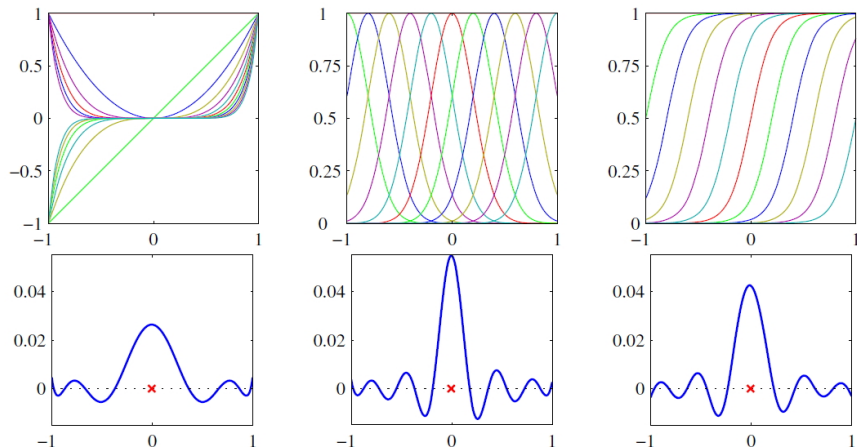
# Constructing kernels

We can define a kernel function by:

- using a nonlinear function $\phi(x)$ to calculate:

$$\kappa(x, x') = \phi(x)^T \phi(x') \tag{11}$$

- using a set of basis functions $\phi_j(x)$, $j = 1, \ldots, M$ and calculate:

$$\kappa(x, x') = \sum_{j=1}^{M} \phi_j(x)^T \phi_j(x') \tag{12}$$

# Constructing kernels



Top: Basis functions using polynomials, Gaussians and logistic sigmoids
Bottom: corresponding $\kappa(x, x')$ plotted as a function of $x$ for $x' = 0$.

# Constructing kernels

Alternatively we can construct a kernel function (expressing similarity) directly, like:

$$\kappa(\mathsf{x}, \mathsf{z}) = \left(\mathsf{x}^T \mathsf{z}\right)^2. \tag{13}$$

However, not any function is a kernel function. It needs to correspond to a dot-product:

$$
\begin{aligned}
\kappa(\mathsf{x}, \mathsf{z}) &= \left(\mathsf{x}^T \mathsf{z}\right)^2 = (x_1 z_1 + x_2 z_2)^2 \\
&= x_1^2 z_1^2 + 2 x_1 z_1 x_2 z_2 + x_2^2 z_2^2 \\
&= (x_1^2, \sqrt{2} x_1 x_2, x_2^2)(z_1^2, \sqrt{2} z_1 z_2, z_2^2) \\
&= \phi(\mathsf{x})^T \phi(\mathsf{z}) \tag{14}
\end{aligned}
$$

A *necessary and sufficient condition* for a kernel function $\kappa(\mathsf{x}, \mathsf{x}')$ to be a valid kernel function is that the corresponding K is positive semi-definite for any choices of $\mathsf{x}_n$, $n = 1, \ldots, N$.

# Constructing kernels

## Techniques for Constructing New Kernels.

Given valid kernels $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$, the following new kernels will also be valid:

$$
\begin{aligned}
k(\mathbf{x}, \mathbf{x}') &= ck_1(\mathbf{x}, \mathbf{x}') \\
k(\mathbf{x}, \mathbf{x}') &= f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \\
k(\mathbf{x}, \mathbf{x}') &= q\left(k_1(\mathbf{x}, \mathbf{x}')\right) \\
k(\mathbf{x}, \mathbf{x}') &= \exp\left(k_1(\mathbf{x}, \mathbf{x}')\right) \\
k(\mathbf{x}, \mathbf{x}') &= k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \\
k(\mathbf{x}, \mathbf{x}') &= k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \\
k(\mathbf{x}, \mathbf{x}') &= k_3\left(\boldsymbol{\phi}(\mathbf{x}), \boldsymbol{\phi}(\mathbf{x}')\right) \\
k(\mathbf{x}, \mathbf{x}') &= \mathbf{x}^{\mathrm{T}}\mathbf{A}\mathbf{x}' \\
k(\mathbf{x}, \mathbf{x}') &= k_a(\mathbf{x}_a, \mathbf{x}_a') + k_b(\mathbf{x}_b, \mathbf{x}_b') \\
k(\mathbf{x}, \mathbf{x}') &= k_a(\mathbf{x}_a, \mathbf{x}_a')k_b(\mathbf{x}_b, \mathbf{x}_b')
\end{aligned}
$$

where $c > 0$ is a constant, $f(\cdot)$ is any function, $q(\cdot)$ is a polynomial with nonnegative coefficients, $\boldsymbol{\phi}(\mathbf{x})$ is a function from $\mathbf{x}$ to $\mathbb{R}^M$, $k_3(\cdot, \cdot)$ is a valid kernel in $\mathbb{R}^M$, $\mathbf{A}$ is a symmetric positive semidefinite matrix, $\mathbf{x}_a$ and $\mathbf{x}_b$ are variables (not necessarily disjoint) with $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$, and $k_a$ and $k_b$ are valid kernel functions over their respective spaces.

# Constructing kernels

Commonly used kernel functions:

- Polynomial kernel of order $M$:

$$\kappa(\mathsf{x}, \mathsf{x}') = \left(\mathsf{x}^T \mathsf{x}' + c\right)^M \tag{15}$$

- Gaussian kernel based on Euclidean distance:

$$\kappa(\mathsf{x}, \mathsf{x}') = exp\left(-\frac{\|\mathsf{x} - \mathsf{x}'\|^2}{2\sigma^2}\right) \tag{16}$$

- Gaussian kernel based on distance function $D(\cdot, \cdot)$:

$$\kappa(\mathsf{x}, \mathsf{x}') = exp\left(-\frac{D(\mathsf{x}, \mathsf{x}')^2}{2\sigma^2}\right) \tag{17}$$

# Kernel Principal Component Analysis

Consider a *centered* data set $\mathcal{D} = \{\phi(x_1), \ldots, \phi(x_N)\}$
(which means $\sum_n \phi(x_n) = 0$).

The data covariance matrix C is:

$$C = \frac{1}{N} \sum_{n=1}^{N} \phi(x_n)\phi(x_n)^T \tag{18}$$

and its eigenvectors are given by solving for $Cv_i = \lambda v_i$.

Using the above:

$$\lambda_i v_i = \frac{1}{N} \sum_{n=1}^{N} \phi(x_n)\Big(\phi(x_n)^T v_i\Big)$$

$$v_i = \sum_{n=1}^{N} \alpha_{in}\phi(x_n) \tag{19}$$

# Kernel Principal Component Analysis

Substituting $v_i$ to C:

$$\frac{1}{N} \sum_{n=1}^{N} \phi(x_n) \phi(x_n)^T \sum_{m=1}^{N} \alpha_{im} \phi(x_m) = \lambda_i \sum_{n=1}^{N} \alpha_{in} \phi(x_n) \qquad (20)$$

Left-multiplying with $\phi(x_l)$ and re-arranging the summations:

$$\frac{1}{N} \sum_{n=1}^{N} \kappa(x_l, x_n) \sum_{m=1}^{N} \alpha_{im} \kappa(x_n, x_m) = \lambda_i \sum_{n=1}^{N} \alpha_{in} \kappa(x_l, x_n) \qquad (21)$$

$$K^2 \boldsymbol{\alpha}_i = \lambda_i N K \boldsymbol{\alpha}_i \qquad (22)$$

Thus, we can calculate $\boldsymbol{\alpha}_i$, $i = 1, \ldots, D'$, by solving the eigenanalysis problem:

$$K \boldsymbol{\alpha}_i = \lambda_i N \boldsymbol{\alpha}_i \qquad (23)$$

# Kernel Principal Component Analysis

Normalization condition:

$$1 = v_i^T v_i = \boldsymbol{\alpha}_i^T K \boldsymbol{\alpha}_i = \lambda_i N \boldsymbol{\alpha}_i^T \boldsymbol{\alpha}_i \tag{24}$$

Thus, we need to normalize $\boldsymbol{\alpha}_i$ with a factor of $\sqrt{\lambda_i N}$.

A data point $x_*$ is mapped to the kernel PCA space by:

$$y_* = A^T k(x_*), \tag{25}$$

where $A = [\boldsymbol{\alpha}_1, \ldots, \boldsymbol{\alpha}'_D] \in \mathbb{R}^{N \times D'}$ and $k(x_*) \in \mathbb{R}^N$ has elements equal to $k_n(x_*) = \kappa(x_n, x_*)$.

# KPCA: Examples



Left: Original data

Right: Transformed data using a $\phi(x)$, and the two principal axes

Kernel PCA with a Gaussian kernel. The contours are lines along which the projection onto the corresponding principal component is constant.

# Centering the kernel matrix K

In order to use a *centered* kernel matrix K, we define:

$$\tilde{\phi}(x_n) = \phi(x_n) - \frac{1}{N} \sum_{l=1}^{N} \phi(x_l) \tag{26}$$

and

$$
\begin{aligned}
\tilde{K}_{nm} &= \tilde{\phi}(x_n)^T \tilde{\phi}(x_m) \\
&= \kappa(x_n, x_m) - \frac{1}{N} \sum_{l=1}^{N} \kappa(x_l, x_m) \\
&\quad - \frac{1}{N} \sum_{l=1}^{N} \kappa(x_n, x_l) + \frac{1}{N^2} \sum_{j=1}^{N} \sum_{l=1}^{N} \kappa(x_j, x_l)
\end{aligned} \tag{27}
$$

Thus:

$$\tilde{K} = K - 1_N K - K 1_N + 1_N K 1_N \tag{28}$$

where $1_N \in \mathbb{R}^{N \times N}$ has elements equal to $1/N$.