# Statistical Learning and Machine Learning
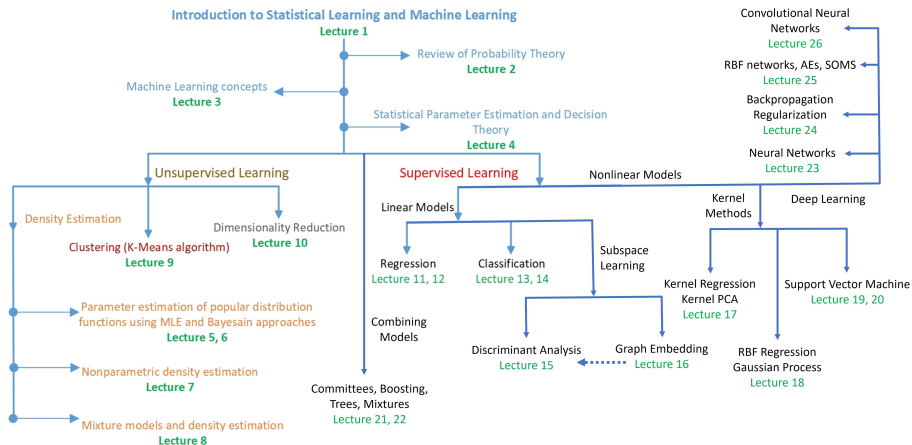## Lecture 18 - Kernel Methods 2

October 13, 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

Deep Learning

Density Estimation

Clustering (K-Means algorithm)
Lecture 9

Dimensionality Reduction
Lecture 10

Linear Models

Kernel Methods

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.

# Constructing kernels
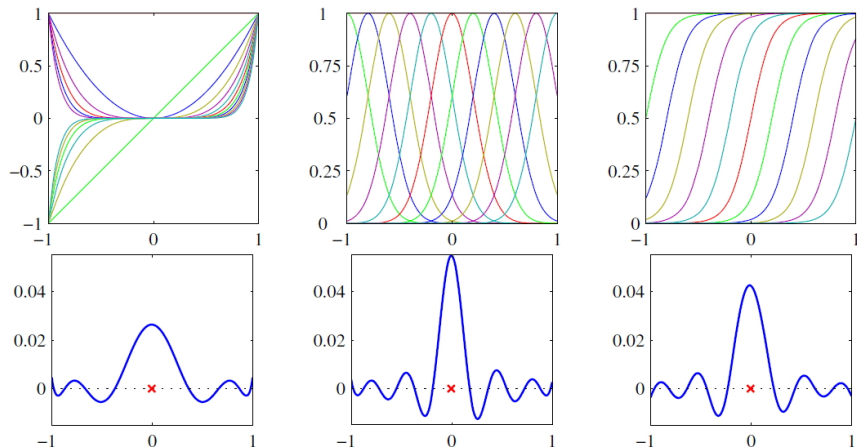
We can define a kernel function by:

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

$$\kappa(x, x') = \phi(x)^T \phi(x') \tag{2}$$

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

$$\kappa(x, x') = \sum_{j=1}^{M} \phi_j(x)^T \phi_j(x') \tag{3}$$

# Constructing kernels



Top: Basis functions using polynomials, Gaussians and ligistic 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:

$$\kappa(\mathsf{x}, \mathsf{z}) = \left(\mathsf{x}^T \mathsf{z}\right)^2. \tag{4}$$

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

$$
\begin{aligned}
\kappa(\mathsf{x}, \mathsf{z}) &= \left(\mathsf{x}^T \mathsf{z}\right)^2 = (x_1 z_1 + x_2 z_2)^2 \\
&= x_1^2 z_1^2 + 2 x_1 z_1 x_2 z_2 + x_2^2 z_2^2 \\
&= (x_1^2, \sqrt{2} x_1 x_2, x_2^2)(z_1^2, \sqrt{2} z_1 z_2, z_2^2) \\
&= \phi(\mathsf{x})^T \phi(\mathsf{z}) \tag{5}
\end{aligned}
$$

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 $M$:

$$\kappa(\mathsf{x}, \mathsf{x}') = \left(\mathsf{x}^T \mathsf{x}' + c\right)^M \tag{6}$$

- Gaussian kernel based on Euclidean distance:

$$\kappa(\mathsf{x}, \mathsf{x}') = exp\left(-\frac{\|\mathsf{x} - \mathsf{x}'\|^2}{2\sigma^2}\right) \tag{7}$$

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

$$\kappa(\mathsf{x}, \mathsf{x}') = exp\left(-\frac{D(\mathsf{x}, \mathsf{x}')^2}{2\sigma^2}\right) \tag{8}$$

# Radial Basis Function Networks

Given a set of data points $x_n \in \mathbb{R}^D$, $n = 1, \ldots, N$ and targets
$t_n \in \mathbb{R}^C$, $n = 1, \ldots, N$, RBF networks apply the following process:

- Select $K$ prototype vectors $\boldsymbol{\mu}_k \in \mathbb{R}^D$, $k = 1, \ldots, K$
- Transform $x_n$, $n = 1, \ldots, N$ to $\boldsymbol{\phi}_n$, $n = 1, \ldots, N \in \mathbb{R}^K$ using a radial basis function $\phi(r)$, where $r = \|x - \boldsymbol{\mu}\|$
- Train a linear regression model from $\tilde{\boldsymbol{\phi}}_n$ to $t_n$:

$$W = \left( \Phi \Phi^T \right)^{-1} \Phi T^T \tag{9}$$

where $\Phi = [\tilde{\boldsymbol{\phi}}_1, \ldots, \tilde{\boldsymbol{\phi}}_N] \in \mathbb{R}^{(K+1) \times N}$, $\tilde{\boldsymbol{\phi}}_n = [1, \boldsymbol{\phi}_n^T]^T$ and $T \in \mathbb{R}^{C \times N}$

# Radial Basis Function Networks

Example Radial Basis Functions

$$\phi(r) = e^{-r^2/2\sigma^2}, \quad \text{Gaussian,}$$

$$\phi(r) = \frac{1}{(\sigma^2 + r^2)^\alpha}, \quad \alpha > 0,$$

$$\phi(r) = \left(\sigma^2 + r^2\right)^\beta, \quad 0 < \beta < 1,$$
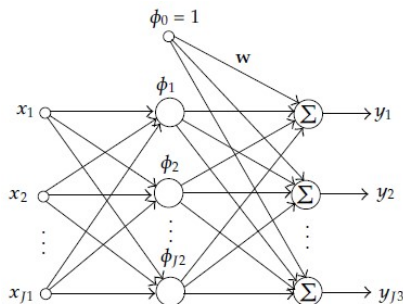
$$\phi(r) = r, \quad \text{linear,}$$

$$\phi(r) = r^2 \ln(r), \quad \text{thin-plate spline,}$$

$$\phi(r) = \frac{1}{1 + e^{(r/\sigma^2)-\theta}}, \quad \text{logistic function,}$$

# Radial Basis Function Networks

The prototype vectors $\boldsymbol{\mu}_k \in \mathbb{R}^D$, $k = 1, \ldots, K$ can be:

- all training data points $x_n$, $n = 1, \ldots, N$. Then $K = N$
- randomly selected training data points $x_n$, $n = 1, \ldots, N$
- $K$ mean vectors obtained by clustering the training data points, e.g. by applying the $K$-Means algorithm

# Gaussian Processes

Let us consider the linear regression model:

$$y(\mathsf{x}) = \mathsf{w}^T \phi(\mathsf{x}) \tag{10}$$

We treat w as a parameter vector drawn from the density function:

$$p(\mathsf{w}) = \mathcal{N}(\mathsf{w}|0, \alpha^{-1}\mathsf{I}). \tag{11}$$

where $\alpha$ is the precision (inverse variance) of the distribution.

Note: The linear regression model is also a function of $\alpha$:

$$y(\alpha, \mathsf{x}) = \mathsf{w}(\alpha)^T \phi(\mathsf{x}) \tag{12}$$

# Gaussian Processes

Given a set of data points $x_n$, $n = 1, \ldots, N$, the joint distribution of the function values is $y = [y(x_1), \ldots, y(x_N)]^T$ and is calculated by:

$$y = \Phi w \tag{13}$$

where $\Phi = [\phi(x_1)^T, \ldots, \phi(x_N)^T]^T$.

Since w is drawn from a Gaussian distribution, y is a Gaussian too:

$$
\begin{aligned}
\mathbb{E}[y] &= \Phi \, \mathbb{E}[w] = 0 \tag{14} \\
cov[y] &= \mathbb{E}[yy^T] = \Phi \, \mathbb{E}[ww^T]\Phi^T = \frac{1}{\alpha}\Phi\Phi^T = K \tag{15}
\end{aligned}
$$

where K is the Gram matrix:

$$K_{nm} = \kappa(x_n, x_m) = \frac{1}{\alpha}\phi(x_n)^T\phi(x_m) \tag{16}$$

and $\kappa(\cdot, \cdot)$ is the kernel function.

# Gaussian Processes: Regression

We model the noise on the observed target values for data point $x_n$ by:

$$t_n = y(x_n) + \epsilon_n \tag{17}$$

where $\epsilon$ is a random noise variable (independent for each data point).

For noise processes following a Gaussian distribution with precision $\beta$:

$$p(t_n|y(x_n)) = \mathcal{N}(t_n|y(x_n), \beta^{-1}) \tag{18}$$

The joint distribution of $t = [t_1, \ldots, t_N]^T$ conditioned on
$y = [y_1, \ldots, y_N]^T$ is:

$$p(t|y) = \mathcal{N}(t|y, \beta^{-1}I_N). \tag{19}$$

# Gaussian Processes: Regression

Using $p(\mathsf{y}) = \mathcal{N}(\mathsf{y}|0, \mathsf{K})$ we get:

$$p(\mathsf{t}) = \int p(\mathsf{t}|\mathsf{y})p(\mathsf{y})d\mathsf{y} = \mathcal{N}(\mathsf{t}|0, \mathsf{C}) \tag{20}$$

where the covariance matrix $\mathsf{C}$ has elements:

$$C(\mathsf{x}_n, \mathsf{x}_m) = \kappa(\mathsf{x}_n, \mathsf{x}_m) + \beta^{-1}\delta_{nm}. \tag{21}$$

$\delta_{nm} = 1$ for $n = m$, and $\delta_{nm} = 0$ for $n \neq m$.

Since the two Gaussian sources of randomness (w.r.t. $y(\mathsf{x})$ and $\epsilon$) are independent, the joint covariance is equal to the summation of the two.

# Gaussian Processes: Regression

Given a new data point $x_{N+1}$, we want to evaluate the $p(t_{N+1}|t)$.

The joint distribution is:

$$p(t_{N+1}) = \mathcal{N}(t_{N+1}|0, C_{N+1}) \tag{22}$$

where

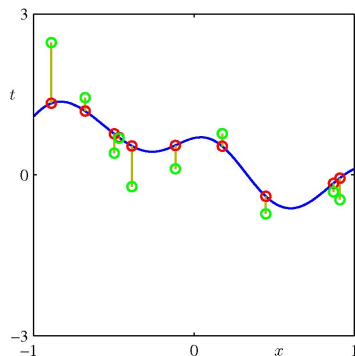$$C_{N+1} = \begin{bmatrix} C_N & k \\ k^T & c \end{bmatrix} \tag{23}$$

with $k = [\kappa(x_1, x_{N+1}), \ldots, \kappa(x_N, x_{N+1})^T]$ and $c = \kappa(x_{N+1}, x_{N+1}) + \beta^{-1}$.

The conditional distribution $p(t_{N+1}|t)$ is a Gaussian with:

$$
\begin{align}
m(x_{N+1}) &= k^T C_N^{-1} t \tag{24} \\
\sigma^2(x_{N+1}) &= c - k^T C_N^{-1} t. \tag{25}
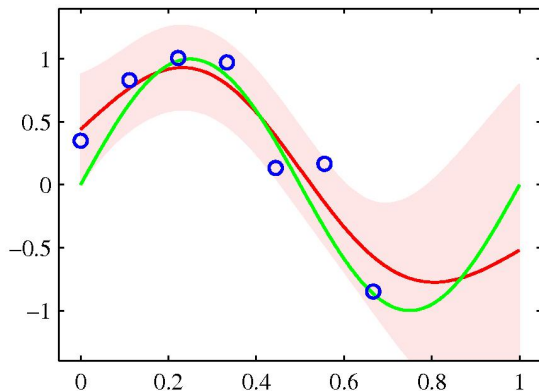\end{align}
$$

# Gaussian Processes: Regression



Line: a sample function from the Gaussian process prior over functions
Red points: the values of $y(x_n)$, $n = 1, \ldots, N$
Green points: targets $t_n$ obtained by adding Gaussian noise to each $y(x_n)$

# Gaussian Processes: Regression



Green line: sinusoidal function generating data points (after adding noise)
Red line: the mean of the Gaussian process predictive distribution
Shaded region: plus and minus two standard deviations on the mean