

# Machine Learning: 708.063 (VO)

## 6. Kernel Methods

May 20, 2020

# Introduction

- ▶ The linear parametric models, we have considered so far can be re-cast into an equivalent **dual** representation.
- ▶ It is based on a linear combination of **kernel functions** evaluated at the training data points.
- ▶ For a fixed nonlinear **feature space** mapping  $\phi(\mathbf{x})$ , the kernel function is given by

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}'),$$

which by definition is a symmetric function in its arguments.

- ▶ The concept of kernels has been introduced into the field of pattern recognition already in 1964 by Aizerman et al.
- ▶ The simple example of a kernel is given by the linear kernel  $\phi(\mathbf{x}) = \mathbf{x}$  such that  $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$ .
- ▶ The concept of defining a kernel by means of an inner product in a feature space is also known as the **kernel trick**.
- ▶ Whenever an algorithm depends on inner products between features, they can be replaced by an equivalent kernel.
- ▶ **Allows to solve problems with infinite-dimensional feature transforms  $\phi(\mathbf{x})$ .**

## Dual representation

- ▶ Many linear models for regression and classification can be reformulated in terms of a dual representation incorporating kernels.
- ▶ Let us consider a linear regression model minimizing a regularized least squares error function

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

- ▶ Can be written in more compact notation as

$$J(\mathbf{w}) = \frac{1}{2} \|\Phi \mathbf{w} - \mathbf{t}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2.$$

- ▶ The dual problem is given by

$$D(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T K \mathbf{a} + \frac{\lambda}{2} \|\mathbf{a} + \mathbf{t}\|^2,$$

where

$$\mathbf{w} = -\frac{1}{\lambda} \Phi^T \mathbf{a}, \quad K = \Phi \Phi^T, \quad \text{with } K_{n,m} = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m).$$

# Discussion

- ▶ The solution of the dual problem is given by

$$\mathbf{a} = -\left(I + \frac{1}{\lambda}K\right)^{-1}\mathbf{t}.$$

- ▶ Substituting back into the regression function gives

$$y(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w} = \mathbf{k}(\mathbf{x})^T (\lambda I + K)^{-1} \mathbf{t},$$

where  $\mathbf{k}(\mathbf{x}) = (k_1(\mathbf{x}), \dots, k_N(\mathbf{x}))^T$  with  $k_n(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}_n)$ .

- ▶ Observe that the regression function is computed without the need to compute the feature transforms  $\phi(\mathbf{x})$ .
- ▶ The original primal problem is a least-squares problem of size  $M \times M$ .
- ▶ The dual problem is again a least squares problem but now of size  $N \times N$ .
- ▶ The problem with the smaller dimension should be considered for solution.

## Constructing kernels

- ▶ One possibility is to compute kernels from the feature representations

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') = \sum_{i=1}^M \phi_i(\mathbf{x}) \phi_i(\mathbf{x}'),$$

- ▶ An alternative approach is to construct kernels directly, but we need to verify that the constructed kernels represent valid kernel functions.
- ▶ As an example, consider a kernel given by

$$k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y})^2.$$

- ▶ In dimension two it can be written as

$$\begin{aligned} k(\mathbf{x}, \mathbf{y}) &= (\mathbf{x}^T \mathbf{y})^2 = (x_1 y_1 + x_2 y_2)^2 = (x_1^2 y_1^2 + 2x_1 y_1 x_2 y_2 + x_2^2 y_2^2 \\ &\quad (x_1^2, \sqrt{2}x_1 x_2, x_2^2)(y_1^2, \sqrt{2}y_1 y_2, y_2^2)^T = \phi(\mathbf{x})^T \phi(\mathbf{y}). \end{aligned}$$

- ▶ Hence, the corresponding feature mapping takes the form

$$\phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1 x_2, x_2^2)^T$$

# Admissible kernels

- ▶ A necessary and sufficient condition for a kernel being admissible is that the corresponding gram matrix with entries  $K_{n,m}$  is positive semi-definite.
- ▶ A standard technique to construct admissible kernels is to consider operations that preserve the validity of a kernel.
- ▶ Let  $k_1(\mathbf{x}, \mathbf{x}')$  and  $k_2(\mathbf{x}, \mathbf{x}')$  be admissible kernels. Then

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

where  $c > 0$ ,  $f$  is any function,  $q$  is a polynomial with nonnegative coefficients,  $\phi(\mathbf{x})$  is a function from  $\mathbb{R}^D$  to  $\mathbb{R}^M$ ,  $k_3(\mathbf{x}, \mathbf{x}')$  is an admissible kernel on  $\mathbb{R}^M$ ,  $A$  is a symmetric and positive definite matrix,  $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$ , and  $k_a, k_b$  are admissible kernels on the partitioned variables.

# Common kernels

- ▶ The kernel  $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^M$  contains only monomials of order  $M$ .
- ▶ The kernel  $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^M$  contains all terms up to order  $M$
- ▶ Another commonly used kernel takes the form

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / (2\sigma^2)),$$

which is called a Gaussian kernel.

- ▶ This is easily shown to be an admissible kernel, as

$$\|\mathbf{x} - \mathbf{x}'\|^2 = \mathbf{x}^T \mathbf{x} + (\mathbf{x}')^T \mathbf{x}' - 2\mathbf{x}^T \mathbf{x}'.$$

- ▶ Hence, the Gaussian kernel can be written as the product

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\mathbf{x}^T \mathbf{x} / (2\sigma^2)) \exp(-(\mathbf{x}')^T \mathbf{x}' / (2\sigma^2)) \exp(\mathbf{x}^T \mathbf{x}' / \sigma^2),$$

which defines a admissible kernel.

# The Nadaraya-Watson model

- Suppose we have training data  $\{\mathbf{x}_n, t_n\}$  and we consider a kernel density estimator to model the joint distribution

$$p(\mathbf{x}, t) = \frac{1}{N} \sum_{n=1}^N f(\mathbf{x} - \mathbf{x}_n, t - t_n),$$

where  $f(\mathbf{x}, t)$  is the kernel function centered at each data point  $(\mathbf{x}_n, t_n)$ .

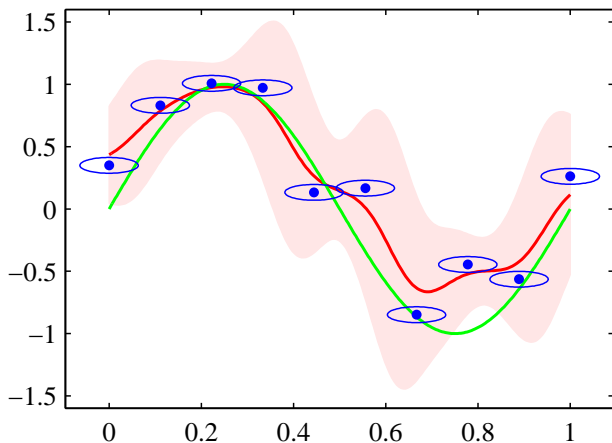
- The regression function  $y(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}]$  can be computed as the conditional mean:

$$y(\mathbf{x}) = \sum_n k(\mathbf{x}, \mathbf{x}_n) t_n, \quad k(\mathbf{x}, \mathbf{x}_n) = \frac{g(\mathbf{x} - \mathbf{x}_n)}{\sum_m g(\mathbf{x} - \mathbf{x}_m)}, \quad g(\mathbf{x}) = \int_{-\infty}^{\infty} f(\mathbf{x}, t) \, dt.$$

- This model is known as the **Nadaraya-Watson** model or **kernel regression**.
- It gives a higher weight to points  $\mathbf{x}_n$  close to the new input  $\mathbf{x}$ .



## Example



- ▶ Illustration of the Nadaraya-Watson kernel regression model.
- ▶ The original sine function is shown in green, and the conditional mean in red.
- ▶ The data points which define the centers of the Gaussians kernels in blue.
- ▶ The shaded area defines the region of plus-minus one standard deviation around the conditional mean.

# Gaussian processes

- ▶ The framework of **Gaussian processes** is obtained by extending the role of kernels to probabilistic discriminative models.
- ▶ We will define a prior probability directly over the functions  $y(\mathbf{x})$ .
- ▶ As a motivation, we start by re-considering the linear regression model

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}),$$

together with a Gaussian prior

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

- ▶ The probability distribution on  $\mathbf{w}$  also induces a probability distribution over the functions  $y(\mathbf{x})$ .
- ▶ Given data  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , we consider the vector-valued function

$$\mathbf{y} = \Phi \mathbf{w},$$

where  $\mathbf{y} = (y(\mathbf{x}_1), \dots, y(\mathbf{x}_N))^T$  and the **design matrix**  $\Phi$  has elements  $\Phi_{nk} = \phi_k(\mathbf{x}_n)$ .

# Joint distribution

- ▶ Since  $\mathbf{w}$  follows a Gaussian distribution, also  $\mathbf{y}$  follows a Gaussian distribution with mean and covariance

$$\mathbb{E}[\mathbf{y}] = \Phi \mathbb{E}[\mathbf{w}] = 0, \quad \text{cov}[\mathbf{y}] = \mathbb{E}[\mathbf{y}\mathbf{y}^T] = \Phi \mathbb{E}[\mathbf{w}\mathbf{w}^T] \Phi^T = \frac{1}{\alpha} \Phi \Phi^T = K,$$

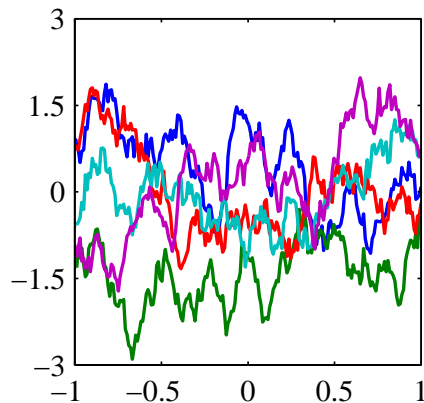
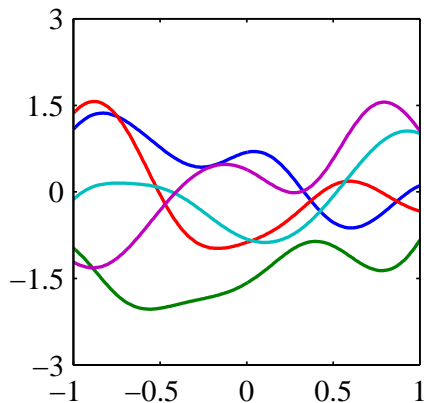
where  $K$  is the Gram matrix with elements

$$K_{n,m} = k(\mathbf{x}_n, \mathbf{x}_m) = \frac{1}{\alpha} \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m),$$

and  $k(\mathbf{x}, \mathbf{x}')$  is the kernel function.

- ▶ We see that a **Gaussian process** is defined as a distribution over functions  $y(\mathbf{x})$  such that the set of functions evaluated at  $\mathbf{x}_1, \dots, \mathbf{x}_N$  are **jointly** Gaussian.
- ▶ Instead of working with features  $\phi(\mathbf{x}_n)$ , one can also directly work with kernel functions.

## Example



- Samples from a Gaussian process with left: a Gaussian kernel and right: an exponential kernel  $k(x, x') = \exp(-\theta|x - x'|)$ .

# Gaussian process for regression

- ▶ We now extend the idea of Gaussian processes to regression under the noise model

$$t_n = y_n + \varepsilon_n.$$

- ▶ We assume that  $\varepsilon_n$  is a Gaussian distributed random variable, therefore

$$p(t_n|y_n) = \mathcal{N}(t_n|y_n, \beta^{-1}).$$

- ▶ The joint distribution of the target vector  $\mathbf{t} = (t_1, \dots, t_N)^T$  conditioned on the values  $\mathbf{y} = (y_1, \dots, y_N)^T$  is therefore given by an isotropic Gaussian of the form

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

where  $I_N$  is the  $N$ -dimensional unit matrix.

- ▶ From the definition of a Gaussian process, the marginal distribution  $p(\mathbf{y})$  is given by a zero-mean Gaussian with covariance defined by the Gram matrix, that is

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|0, K),$$

where  $K$  is again the Gram matrix.

# The marginal distribution

- ▶ The marginal distribution is obtained by integrating the joint distribution  $p(\mathbf{y}, \mathbf{t}) = p(\mathbf{t}|\mathbf{y})p(\mathbf{y})$  with respect to  $\mathbf{y}$ :

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

where the covariance matrix  $\mathbf{C}$  has elements

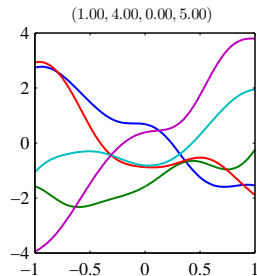
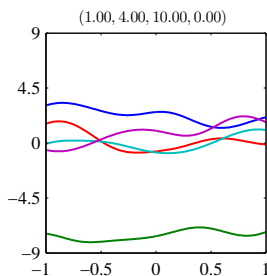
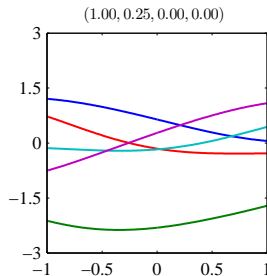
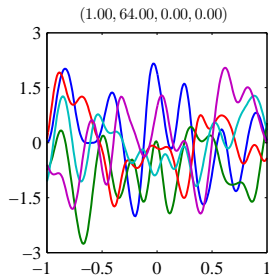
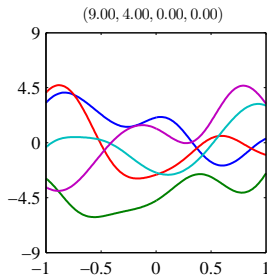
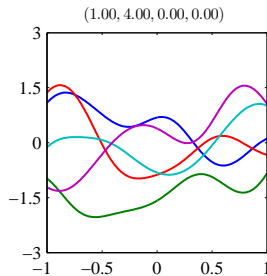
$$\mathbf{C}(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \beta^{-1}\delta_{nm}.$$

- ▶ This covariance reflects the fact that the Gaussian nature of the functions  $\mathbf{y}(\mathbf{x})$  are independent from the Gaussian noise assumption on  $\varepsilon$ .
- ▶ A widely used kernel function for Gaussian process regression is given by

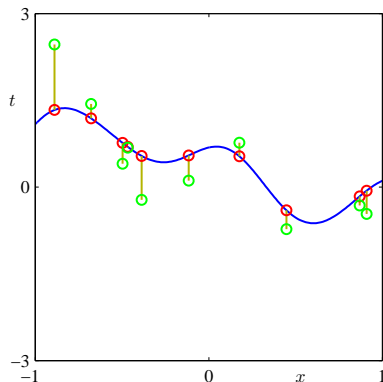
$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{\theta_1}{2}\|\mathbf{x}_n - \mathbf{x}_m\|^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m,$$

which is a combination of a Gaussian kernel, a constant term and a linear term.

# Examples of Gaussian process priors



# Sampling from a Gaussian process



- ▶ The blue curve shows a sample function from the Gaussian process prior  $p(\mathbf{y})$ .
- ▶ The red points are obtained by evaluating the sampled function values  $y_n$  on a set of input vectors  $x_n$ .
- ▶ The green points are obtained by adding independent Gaussian noise to the  $y_n$ .



# Making predictions

- ▶ Based on a training data set  $\mathbf{t}_N = (t_1, \dots, t_N)^T$  and corresponding input vectors  $(\mathbf{x}_1, \dots, \mathbf{x}_N)$  the goal is to predict the target variable  $t_{N+1}$  for a new input vector  $\mathbf{x}_{N+1}$ .
- ▶ For this, we need to evaluate the predictive distribution

$$p(t_{N+1}|\mathbf{t})$$

where the dependence on the  $\mathbf{x}_n$  is omitted.

- ▶ We start by writing down the marginal distribution

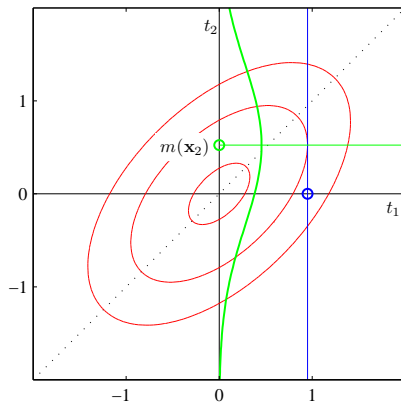
$$p(\mathbf{t}_{N+1}) = \mathcal{N}(\mathbf{t}_{N+1}|0, C_{N+1}), \quad C_{N+1} = \begin{pmatrix} C_N \mathbf{k} \\ \mathbf{k}^T c \end{pmatrix},$$

where  $\mathbf{k} = (k(\mathbf{x}_1, \mathbf{x}_{N+1}), \dots, k(\mathbf{x}_N, \mathbf{x}_{N+1}))^T$  and  $c = k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1}$ .

- ▶ The predictive distribution  $p(t_{N+1}|\mathbf{t})$  is a conditional Gaussian and is computed as

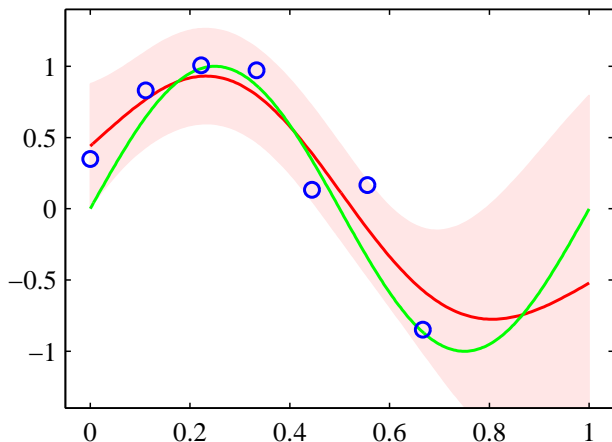
$$p(t_{N+1}|\mathbf{t}) = \mathcal{N}(t_{N+1} | \underbrace{\mathbf{k}^T C_N^{-1} \mathbf{t}}_{m(\mathbf{x}_{N+1})}, \underbrace{c - \mathbf{k}^T C_N^{-1} \mathbf{k}}_{\sigma^2(\mathbf{x}_{N+1})})$$

## Example



- ▶ The horizontal axis is  $t_1$ , the vertical axis is  $t_2$ .
- ▶ The red ellipses are the contours of the joint distribution  $p(t_1, t_2)$ .
- ▶ The blue line  $t_1$  is the training data point
- ▶ The conditional distribution  $p(t_2 | t_1)$  is shown as the green curve.

## Gaussian process regression



- ▶ The green curve is the true sinusoidal function.
- ▶ The blue points are the noisy data points  $(x_n, t_n)$ .
- ▶ The red line is the mean of the Gaussian process.
- ▶ The shaded regions corresponds to plus-minus one standard deviation around the mean.