

# Machine Learning for Physics and Astronomy

**Juan Rojo**

VU Amsterdam & Theory group, Nikhef

***Natuur- en Sterrenkunde BSc (Joint Degree), Honours Track***

***Lecture 7, 12/10/2020***

# Today's lecture

- 📌 Kernel methods and the Dual Representation
- 📌 Gaussian processes
- 📌 Guest lecture by **Dr. Tristan Bereau** on applications of machine learning in condensed matter

# Kernel Methods

# Kernel Methods

most of the ML models we have considered for regression and classification are based on the **parametric mapping** (linear or non-linear) between inputs and outputs

$$y(x, \theta)$$

were the **training examples are discarded** after the model parameters (or their posterior distribution) have been determined

here we introduce models where the training dataset is also used in the **prediction phase**

recall discussion of **Support Vector Machines**: we introduced a **feature space mapping**

*Example of non-linear feature-space transformation*

$x$	$\phi(x)$
<i>Height</i>	<i>Height/Weight</i>
<i>Weight</i>	<i>Weight+Height</i>
<i>Age</i>	<i>Age * Weight</i>

# Kernel Methods

most of the ML models we have considered for regression and classification are based on the **parametric mapping** (linear or non-linear) between inputs and outputs

$$y(x, \theta)$$

were the **training examples are discarded** after the model parameters (or their posterior distribution) have been determined

here we introduce models where the training dataset is also used in the **prediction phase**

the key ingredient of **kernel methods** is the **kernel function** evaluated at the training points

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

*kernel function* *feature space mapping*

the simplest kernel is the **linear kernel** (identity mapping in feature space)

$$k(x, x') = x^T x'$$

# Dual representations

many ML models for regression and classification can be reformulated in a **dual representation** where the kernel functions arise naturally

consider a linear regression model with a regularised least-squares error function

$$E_{\text{tr}}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{n=1}^N (\boldsymbol{\theta}^T \boldsymbol{\phi}(x_n) - t_n)^2 + \frac{\lambda}{2} \boldsymbol{\theta}^T \boldsymbol{\theta}$$

*regulator* (points to  $\frac{\lambda}{2} \boldsymbol{\theta}^T \boldsymbol{\theta}$ )

*output from training examples* (points to  $t_n$ )

the model parameters are determined **analytically** by requiring the vanishing of the gradient

$$\boldsymbol{\theta} = \sum_{n=1}^N a_n \boldsymbol{\phi}(x_n), \quad a_n = -\frac{1}{\lambda} (\boldsymbol{\theta}^T \boldsymbol{\phi}(x_n) - t_n)$$

*recall that inner products take place in feature space!*

one can formulate a dual representation of the problem in terms of the parameter vector

$$\boldsymbol{a} = (a_1, a_2, \dots, a_N)^T \quad \boldsymbol{t} = (t_1, t_2, \dots, t_N)^T$$

# Dual representations

with some algebra one can show that the original figure of merit of the model

$$E_{\text{tr}}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{n=1}^N (\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}_n) - t_n)^2 + \frac{\lambda}{2} \boldsymbol{\theta}^T \boldsymbol{\theta}$$

admits a **dual representation** in terms of **kernel function**

$$E_{\text{tr}}(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}$$

where  $\mathbf{K}$  is known as the Gram matrix, an  $N \times N$  symmetric matrix with entries

$$K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m) = \boldsymbol{\phi}(\mathbf{x}_n)^T \boldsymbol{\phi}(\mathbf{x}_m)$$

given by the kernel function evaluated over two of the **input training examples**

the crucial property of this dual representation is that now the predictions for new inputs will explicitly **depend on the training examples**

# Dual representations

the crucial property of this dual representation is that now the predictions for new inputs will explicitly **depend on the training examples**

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

*output of trained model* →  $y(\mathbf{x})$

*new input* →  $\mathbf{x}$

$$\mathbf{k}(\mathbf{x}) = (k(\mathbf{x}_1, \mathbf{x}), k(\mathbf{x}_2, \mathbf{x}), \dots)$$

in the dual formulation we invert a  $N \times N$  matrix (data space) rather than a  $M \times M$  one (feature space), so this does not appear to be advantageous ...

the main benefit is being able to **work directly with kernel functions** and bypass a choice of feature map

this allows one to work in feature spaces of **very high (even infinite) dimensionality**



# Choosing kernel functions

recall that the kernel function is constructed from a **feature space mapping**

$$k(\mathbf{x}_n, \mathbf{x}_m) = \boldsymbol{\phi}(\mathbf{x}_n)^T \boldsymbol{\phi}(\mathbf{x}_m) = \sum_{i=1}^M \phi_i(\mathbf{x}_n) \phi_i(\mathbf{x}_m)$$

one can also construct kernel functions directly, provided they can be expressed as above

*e.g., for a 2D input space*

$$k(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2 = (x_1 z_1 + x_2 z_2)^2 = \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{z})$$

$$\boldsymbol{\phi}(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

*Note that feature space dimension (M=3) is bigger than input space dimensionality (d=2)!*

there exist more general methods to construct acceptable kernels e.g. using **basis functions**

a popular choice is Gaussian kernel, built from an **infinite-dimensional feature map**

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

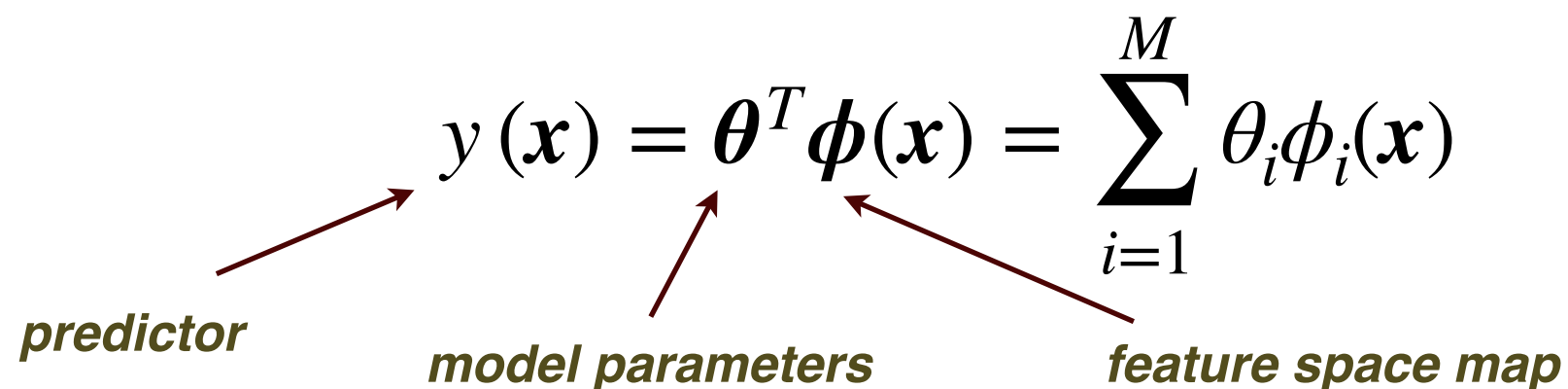
# Gaussian Processes

# Gaussian processes

the previous discussion was limited to **non-probabilistic models**. We now extend the kernel methods to **probabilistic discrimination processes**

in Gaussian processes we bypass the need of a **parametric classifier/regressor** and work directly in the **space of functions for the (prior) probability distributions**

we can illustrate the philosophy of the method by revisiting the simple linear regression example

$$y(\mathbf{x}) = \boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}) = \sum_{i=1}^M \theta_i \phi_i(\mathbf{x})$$


*predictor*                      *model parameters*                      *feature space map*

assume that the model parameters have as **prior distribution** an isotropic Gaussian

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} \mid \mathbf{0}, \alpha^{-1} \mathbf{1})$$

for a given choice of the model parameters we will have a different functional form of the predictor

# Gaussian processes