

KERNEL METHODS

Chia-Ping Chen

Professor

National Sun Yat-sen University

Department of Computer Science and Engineering

Machine Learning

- Dual Representation
- Constructing Kernels
- Radial Basis Functions
- Gaussian Processes
- Maximum Margin Classifiers
- Relevance Vector Machines

With a kernel method, the prediction for an unseen \boldsymbol{x} is

$$y(\boldsymbol{x}) = \sum_n a_n k(\boldsymbol{x}_n, \boldsymbol{x})$$

- $k(\boldsymbol{x}, \boldsymbol{x}')$ is a **kernel function**
- n runs through a data set
- a_n are trainable parameters

- Gaussian kernel

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

- exponential kernel

$$k(x, x') = \exp(-\theta|x - x'|)$$

Consider a linear regression model $y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$. Let the model parameters \mathbf{w} be trained with $\mathcal{D} = \{(\mathbf{x}_n, t_n)\}_{n=1}^N$ by minimizing

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

where $\lambda > 0$. The solution for \mathbf{w} takes the form

$$\mathbf{w} = -\frac{1}{\lambda} \sum_{n=1}^N [\mathbf{w}^T \phi(\mathbf{x}_n) - t_n] \phi(\mathbf{x}_n) = \Phi^T \mathbf{a}$$

where Φ is the design matrix and \mathbf{a} has components

$$a_n = -\frac{1}{\lambda} [\mathbf{w}^T \phi(\mathbf{x}_n) - t_n]$$

With $\mathbf{w} = \Phi^T \mathbf{a}$, the prediction function of linear regression is equivalent to a kernel machine

$$\begin{aligned}y(\mathbf{x}) &= \mathbf{w}^T \phi(\mathbf{x}) \\&= \mathbf{a}^T \Phi \phi(\mathbf{x}) \\&= \sum_{n=1}^N a_n \phi(\mathbf{x}_n)^T \phi(\mathbf{x}) \\&= \sum_{n=1}^N a_n k(\mathbf{x}_n, \mathbf{x})\end{aligned}$$

with the kernel function related to the basis functions by

$$k(\mathbf{x}_n, \mathbf{x}) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x})$$

If we substitute $\mathbf{w} = \Phi^T \mathbf{a}$, the objective function $J(\mathbf{w})$ becomes

$$\begin{aligned} J(\mathbf{a}) &= \frac{1}{2} \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a} \\ &= \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a} \end{aligned}$$

where $\mathbf{K} = \Phi \Phi^T$ is the Gram matrix. Note \mathbf{K} is positive definite and

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

Setting the gradient of $J(\mathbf{a})$ to zero, we obtain

$$\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

The prediction is

$$y(\mathbf{x}) = \sum_{n=1}^N a_n k(\mathbf{x}_n, \mathbf{x}) = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

where $\mathbf{k}(\mathbf{x})$ has components

$$k_n(\mathbf{x}) = k(\mathbf{x}_n, \mathbf{x})$$

We invert an $N \times N$ matrix in determining \mathbf{a} of a kernel machine

$$y(\mathbf{x}) = \sum_{n=1}^N a_n k(\mathbf{x}_n, \mathbf{x})$$

In comparison, we invert an $M \times M$ matrix in determining \mathbf{w} of a linear regression model

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

Constructing Kernel Functions

In dual representation of linear regression, we have defined a kernel function as the inner product of the feature vectors

$$k(\mathbf{x}_n, \mathbf{x}) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x})$$

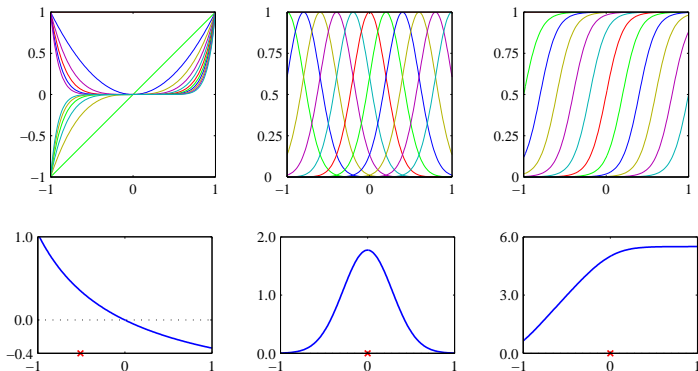
Thus, one way to define a kernel function is use the inner product of a feature space.

KERNEL AND BASIS

For one-dimensional input space, we have

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

where $\phi_i(x)$ are the basis functions.



Instead of defining kernel functions based on basis functions, we can also define kernel functions directly.

For example, we can define

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^2$$

This is equivalent to a particular set of basis functions, since

$$(\mathbf{x}^T \mathbf{x}')^2 = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

where

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

The point is that the feature space is **implicit** and not represented.

KERNEL CONSTRUCTION

Kernels can be constructed from kernels in various ways.

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}'))$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}'))$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}'), \phi(\mathbf{x}'))$$

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

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b)$$

where $c > 0$, $f(\cdot)$ is any function, $q(\cdot)$ is polynomial with non-negative coefficients, $\phi(\cdot)$ is feature function, and $(\mathbf{x}_a, \mathbf{x}_b) = \mathbf{x}$.

Radial Basis Functions

A radial function is a function whose value at a point depends on the radial distance of the point from a center. That is

$$\phi(\mathbf{x}) = h(\|\mathbf{x} - \boldsymbol{\mu}\|)$$

where $\boldsymbol{\mu}$ is the center.

LINEAR REGRESSION WITH RADIAL BASIS FUNCTION

Suppose we have input values $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ along with target values $\{t_1, \dots, t_N\}$. For linear regression, we can define N basis functions each centered at an input point

$$f(\mathbf{x}) = \sum_{n=1}^N w_n h(\|\mathbf{x} - \mathbf{x}_n\|)$$

Then $\{w_1, \dots, w_N\}$ can be found to fit every data point exactly, i.e.

$$t_i = \sum_{n=1}^N w_n h(\|\mathbf{x}_i - \mathbf{x}_n\|)$$

Suppose we have input values $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ along with target values $\{t_1, \dots, t_N\}$. A special form of prediction based on kernel functions is

$$y(\mathbf{x}) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n)t_n$$

In this form, the kernel function $k(\cdot, \cdot)$ has parameters which are decided by data.

Let $\mathcal{D} = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\}$ be a data set. Let the joint probability $p(\mathbf{x}, t)$ be estimated by

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

where $f(\mathbf{x}, t)$ is a component density function with

$$\iint f(\mathbf{x}, t) d\mathbf{x} dt = 1$$

and

$$\int t f(\mathbf{x}, t) dt = 0$$

OPTIMUM REGRESSION

The prediction of t given \mathbf{x} with the minimum expected squared error is the conditional mean

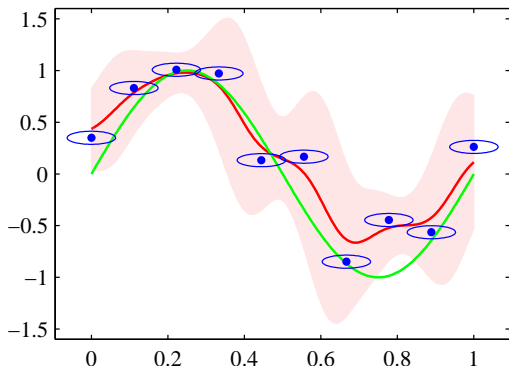
$$\begin{aligned}y(\mathbf{x}) &= \mathbb{E}[t|\mathbf{x}] = \int tp(t|\mathbf{x})dt = \frac{\int tp(\mathbf{x}, t)dt}{\int p(\mathbf{x}, t)dt} \\&= \frac{\sum_n \int tf(\mathbf{x} - \mathbf{x}_n, t - t_n)dt}{\sum_m \int f(\mathbf{x} - \mathbf{x}_m, t - t_m)dt} \\&= \frac{\sum_n g(\mathbf{x} - \mathbf{x}_n)t_n}{\sum_m g(\mathbf{x} - \mathbf{x}_m)}, \quad g(\mathbf{x}) = \int f(\mathbf{x}, t)dt \\&= \sum_n k(\mathbf{x}, \mathbf{x}_n)t_n\end{aligned}$$

where

$$k(\mathbf{x}, \mathbf{x}_n) = \frac{g(\mathbf{x} - \mathbf{x}_n)}{\sum_m g(\mathbf{x} - \mathbf{x}_m)}$$

DISCUSSION AND EXAMPLE

- Nadaraya-Watson model leads to a kernel regression
- The kernel function is a radial function
- Example: isotropic Gaussian component density



Gaussian Processes

- A **stochastic process** or **random process** is a collection of indexed random variables.
- A stochastic process is characterized by the joint distribution of any finite set of random variables in the process.
- A **Gaussian process** is a stochastic process such that the joint distribution of any finite set of random variables is Gaussian.

A Gaussian process is completely specified by the means and covariances of the random variables in the process.

- This is because a Gaussian distribution is specified by a mean vector and a covariance matrix.
- Let $y(t)$ be a Gaussian process. Then it is specified by

$$\mu(t) = \mathbb{E}[y(t)]$$

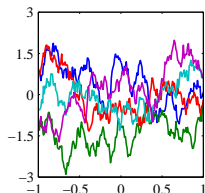
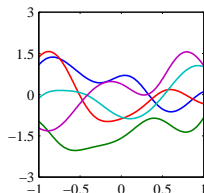
and

$$k(t, t') = \text{cov}(y(t), y(t'))$$

EXAMPLES OF GAUSSIAN PROCESS

A zero-mean Gaussian process is specified by a **kernel function** for the covariance

$$k(t, t') = \mathbb{E}[y(t)y(t')]$$



Gaussian processes with Gaussian (left) and exponential kernel.

A linear regression function in which the parameters have a Gaussian distribution is a Gaussian process.

Consider $y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$ with $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{0}, \alpha^{-1} \mathbf{I})$.

- For any $\mathbf{y} = (y(\mathbf{x}_1), \dots, y(\mathbf{x}_N))^T$, we have $\mathbf{y} = \Phi \mathbf{w}$. As \mathbf{y} is a linear function of \mathbf{w} , it is Gaussian.
- Thus $y(\mathbf{x})$ is a Gaussian process.
- The mean and kernel (covariance) functions of $y(\mathbf{x})$ are

$$\mu(\mathbf{x}) = \mathbb{E}[y(\mathbf{x})] = 0$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[y(\mathbf{x})y(\mathbf{x}')] = \alpha^{-1} \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

How do we fit a data set to a Gaussian process?

Let $\mathcal{D} = \{(\mathbf{x}_n, t_n)\}_{n=1}^N$ be a data set.

- We assume a zero-mean Gaussian process $y(\mathbf{x})$ with kernel function $k(\mathbf{x}, \mathbf{x}')$.
- We assume noises on the target values $t_n = y(\mathbf{x}_n) + \epsilon_n$, where ϵ_n are i.i.d. Gaussian noises with variance β^{-1} . That is

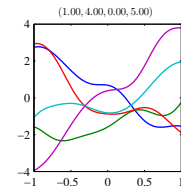
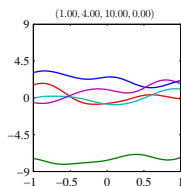
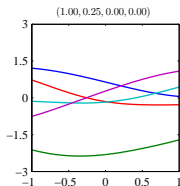
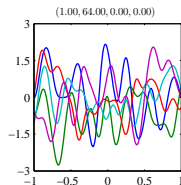
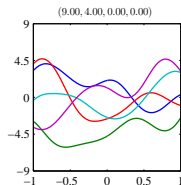
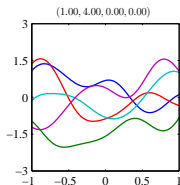
$$p(t_n|y_n) = \mathcal{N}(t_n|y_n, \beta^{-1}), \quad y_n = y(\mathbf{x}_n)$$

- The parameters in the kernel function and the precision β can be learned from \mathcal{D} .

PARAMETRIC KERNEL FUNCTION

One widely used parametric kernel function is

$$k(\mathbf{x}, \mathbf{x}') = \theta_0 \exp \left\{ -\frac{\theta_1}{2} \|\mathbf{x} - \mathbf{x}'\|^2 \right\} + \theta_2 + \theta_3 \mathbf{x}^T \mathbf{x}'$$



PROBABILITY OF A VECTOR OF TARGETS

Denote $\mathbf{y} = [y_1, \dots, y_N]^T$ and $\mathbf{t} = [t_1, \dots, t_N]^T$. We have

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

and

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

Thus

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

where

$$\mathbf{C} = \mathbf{K} + \beta^{-1}\mathbf{I}_N$$

Given a Gaussian process $y(\mathbf{x})$ and the parameter β of Gaussian noise on targets, how do we predict t for input \mathbf{x} ?

- The joint probability of $\mathbf{t}_{N+1} = [\mathbf{t}_N, t]^T$ is

$$p(\mathbf{t}_{N+1}) = \mathcal{N}(\mathbf{t}_{N+1}; \mathbf{0}, \mathbf{C}_{N+1})$$

- Note $\mathbf{C}_{N+1} = \begin{pmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & c \end{pmatrix}$ where $\mathbf{k} = [k_1, \dots, k_N]^T$ with $k_n = k(\mathbf{x}_n, \mathbf{x})$ and $c = k(\mathbf{x}, \mathbf{x}) + \beta^{-1}$.
- It follows that $p(t|\mathbf{t}_N) = \mathcal{N}(t; m(\mathbf{x}), \sigma^2(\mathbf{x}))$, where $m(\mathbf{x}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}_N$ and $\sigma^2(\mathbf{x}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}$.
- Thus, the optimal prediction is a kernel regression function

$$m(\mathbf{x}) = \sum_{n=1}^N a_n k(\mathbf{x}_n, \mathbf{x}), \quad \mathbf{a} = \mathbf{C}_N^{-1} \mathbf{t}_N$$