

# GAUSSIAN PROCESSES

2009年11月26日

上午 11:02

$$y(x) = W^T \phi(x) \quad \text{prior over } W: p(W) = \mathcal{N}(W | 0, \alpha^{-1} I)$$

↑  
precision

probability distribution over  $W$   
induces probability distribution over  $y(x)$

In practice, we evaluate this function at specific values of  $x$ :  $x_1, \dots, x_N$ , we get  $y(x_1), \dots, y(x_N)$

What does the joint distribution of  $y(x_1), \dots, y(x_N)$  look like?

$$y = \Phi W \quad \Phi_{nk} = \phi_k(x_n)$$

linear combination of Gaussian distributed variables

so  $y$  is Gaussian

$$E[y] = \Phi E[W] = 0$$

$$\text{cov}[y] = E[(y-0)(y-0)^T] = \Phi E[WW^T] \Phi^T = \frac{1}{\alpha} \Phi \Phi^T = K$$

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

① A Gaussian process is defined as a probability distribution over functions  $y(x)$  such that the set of values of  $y(x)$  evaluated at an arbitrary set of point  $x_1, \dots, x_N$  jointly have a Gaussian distribution.

② The joint distribution over  $N$  variables  $y_1, \dots, y_N$  is specified completely by the mean and covariance

mean function and covariance function

2009年11月26日

上午 11:51

# Gaussian Processes for Regression

$$t_n = y_n + \epsilon_n \quad y_n = y(x_n) \quad \epsilon_n \text{ is a random noise variable}$$

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

↑  
precision of the noise

$$p(t | y) = \mathcal{N}(t | y, \beta^{-1} I_N) \quad \text{isotropic Gaussian}$$

noise is independent for each data point

$$p(y) = \mathcal{N}(y | 0, K) \quad (\text{from the previous pages})$$

we want to find the marginal

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

$$\Rightarrow C = \beta^{-1} I_N + K$$

$$\left. \begin{aligned} p(a) &= \mathcal{N}(a | \mu, \Lambda^{-1}) \\ p(b|a) &= \mathcal{N}(b | Aa + d, L^{-1}) \\ p(b) &= \mathcal{N}(b | A\mu + d, L^{-1} + AL^{-1}A^T) \end{aligned} \right\}$$

2009年11月26日  
下午 10:02

make predictions

$$t_N = (t_1, \dots, t_N)^T$$

predictive distribution  $p(t_{N+1} | t_N)$

start by writing down the joint distribution  
 $p(t_{N+1}) \rightarrow p(t_{N+1} | t_N)$

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

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

$$c = k(x_{N+1}, x_{N+1}) + \beta^{-1}$$

$$k = \begin{pmatrix} \vdots \\ k(x_n, x_{N+1}) \\ \vdots \end{pmatrix}$$

$$\mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} = \begin{pmatrix} t \\ t_{N+1} \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$$

$$\mu_{b|a} = \mu_b + \Sigma_{ba} \Sigma_{aa}^{-1} (x_a - \mu_a)$$

$$\Sigma_{b|a} = \Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab}$$

Joint  $\rightarrow$  conditional

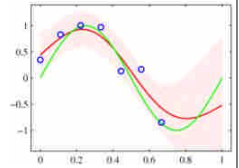
$$m(x_{N+1}) = k^T C_N^{-1} t$$

$$\sigma^2(x_{N+1}) = c - k^T C_N^{-1} k$$

2009年11月26日  
下午 10:10

The predictive distribution is a Gaussian  
whose mean and variance both depend on  $x_{N+1}$

(Fig. 6.8)



$\lambda_i$  eigenvalue of  $K$   
then  $\lambda_i + \beta^{-1}$  eigenvalue of  $C$

$$m(x_{N+1}) = \sum_{n=1}^N a_n k(x_n, x_{N+1})$$

$a_n$  is the  $n$ th component of  $C_N^{-1} t$

e.g.  $k(x, x') = \phi(\|x - x'\|)$  RBF

for kernel methods,  $C$  is an  $N \times N$  matrix  
inversion needs  $O(N^3)$

for linear regression  $S_N$  is  $M \times M$   
inversion needs  $O(M^3)$   $S_N^{-1} = \alpha I + \beta \Phi^T \Phi$

what if the data set is large?

sparse Gaussian processes

# Gaussian Processes for Classification

2009年11月26日

下午 10:20

$t \in \{0, 1\}$   $y = \sigma(a)$   $a(x)$  is linear

$$p(t|a) = \sigma(a)^t (1 - \sigma(a))^{1-t}$$

$$\begin{cases} p(a_{N+1}) = \mathcal{N}(a_{N+1} | 0, C_{N+1}) \\ C(x_n, x_m) = k(x_n, x_m) + \nu \delta_{nm} \end{cases} \quad \delta_{nm} = \begin{cases} 1 & n=m \\ 0 & n \neq m \end{cases}$$

we want to get  $p(t_{N+1}=1 | t_N) = \int p(t_{N+1}=1 | a_{N+1}) \underbrace{p(a_{N+1} | t_N)}_{\sigma(a_{N+1})} da_{N+1}$

this integration is intractable

we may apply Laplace approximation  
first, for the second term

$$\begin{aligned} p(a_{N+1} | t_N) &= \int p(a_{N+1}, \bar{a}_N | t_N) d\bar{a}_N \\ &= \frac{1}{p(t_N)} \int p(a_{N+1}, \bar{a}_N) p(t_N | a_{N+1}, \bar{a}_N) d\bar{a}_N \\ &= \frac{1}{p(t_N)} \int p(a_{N+1} | \bar{a}_N) p(\bar{a}_N) p(t_N | \bar{a}_N) d\bar{a}_N \quad \leftarrow \text{use independence} \\ &= \int p(a_{N+1} | \bar{a}_N) p(\bar{a}_N | t_N) d\bar{a}_N \end{aligned}$$

from GP regression, we know  $p(a_{N+1} | \bar{a}_N) = \mathcal{N}(a_{N+1} | k_{CN}^T \bar{a}_N, c - k_{CN}^T C_N^{-1} k)$

so we need to find a Gaussian to fit  $p(\bar{a}_N | t_N)$

2009年11月26日

下午 10:46

$p(\bar{a}_N)$  is zero-mean Gaussian

$$\begin{aligned} p(t_N | \bar{a}_N) &= \prod_{n=1}^N \sigma(a_n)^{t_n} (1 - \sigma(a_n))^{1-t_n} \\ &= \prod_{n=1}^N \left( \frac{1}{1+e^{-a_n}} \right)^{t_n} \left( \frac{e^{-a_n}}{1+e^{-a_n}} \right)^{1-t_n} \\ &= \prod_{n=1}^N \left( \frac{1}{1+e^{-a_n}} \right)^{t_n} \left( \frac{1}{1+e^{-a_n}} \right)^{1-t_n} (e^{-a_n})^{1-t_n} \\ &= \prod_{n=1}^N \frac{1}{1+e^{-a_n}} e^{-a_n} \cdot e^{a_n t_n} \\ &= \prod_{n=1}^N e^{a_n t_n} \frac{1}{1+e^{a_n}} = \prod_{n=1}^N e^{a_n t_n} \sigma(-a_n) \end{aligned}$$

∅ find the mode

$$\begin{aligned} \Psi(\bar{a}_N) &= \ln p(\bar{a}_N | t_N) = \ln p(\bar{a}_N) + \ln p(t_N | \bar{a}_N) \\ &= -\frac{1}{2} \bar{a}_N^T C_N^{-1} \bar{a}_N - \frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln |C_N| + t_N^T \bar{a}_N \\ &\quad - \sum_{n=1}^N \ln(1+e^{a_n}) \end{aligned}$$

mode  $\nabla \Psi(\bar{a}_N) = t_N - \sigma_N - C_N^{-1} \bar{a}_N$

$$\sigma_N = \begin{pmatrix} \vdots \\ \sigma(a_n) \\ \vdots \end{pmatrix}$$

optimized by iterative reweighted least squares  
IRLS

2009年11月26日  
下午 11:01

compute the Hessian for IRLS

$$\nabla \nabla \Psi(\bar{a}_N) = -W_N - C_N^{-1}$$

$$W_N = \begin{bmatrix} \ddots & & \\ \sigma(a_n)(1-\sigma(a_n)) & & \\ \vdots & & \ddots \end{bmatrix} \text{ positive definite}$$

$\Rightarrow -\nabla \nabla \Psi(\bar{a}_N)$  positive definite

$p(\bar{a}_N | t_N)$  is log convex  $\Rightarrow$  single mode, global maximum

$$\bar{a}_N^{\text{new}} = C_N (I + W_N C_N)^{-1} \{ t_N - \sigma_N + W_N \bar{a}_N \}$$

$$\bar{a}^* = C_N (t_N - \sigma_N) \mid \nabla \Psi(\bar{a}_N) = 0$$

$$\text{Hessian } H = -\nabla \nabla \Psi(\bar{a}_N^*) = W_N + C_N^{-1} \text{ evaluated at } \bar{a}^*$$

$$\text{so } q(\bar{a}_N) = \mathcal{N}(\bar{a}_N \mid \bar{a}_N^*, H_{\bar{a}^*}^{-1})$$

is an approximation to  $p(\bar{a}_N | t_N)$

$$\text{back to } p(a_{N+1} | t_N) = \int p(a_{N+1} | \bar{a}_N) p(\bar{a}_N | t_N) d\bar{a}_N$$

$$\simeq \int \underset{\text{Gaussian}}{p(a_{N+1} | \bar{a}_N)} \underset{\text{Gaussian}}{q(\bar{a}_N)} d\bar{a}_N$$

$$\text{recall } p(a_{N+1} | \bar{a}_N) = \mathcal{N}(a_{N+1} \mid k^T C_N^{-1} \bar{a}_N, c - k^T C_N^{-1} k)$$

combined with  $q(\bar{a}_N)$

$$\text{we get } p(a_{N+1} | t_N) \simeq \mathcal{N}(\mu, s^2)$$

where

$$\mu = k^T (t_N - \sigma_N), \quad s^2 = c - k^T (W_N^{-1} + C_N)^{-1} k$$

2009年11月26日  
下午 11:24

Finally, the predictive distribution

$$\begin{aligned} p(t_{N+1} = 1 \mid t_N) &= \int p(t_{N+1} = 1 \mid a_{N+1}) p(a_{N+1} \mid t_N) da_{N+1} \\ &\quad \downarrow \qquad \qquad \qquad \downarrow \\ &\quad \sigma(a_{N+1}) \qquad \qquad \mathcal{N}(\mu, s^2) \\ &\simeq \sigma(k(\sigma^2)\mu) \qquad (4.153) \end{aligned}$$

We skip the part of deciding parameters for covariance function

# SPARSE KERNEL MACHINES

2009年11月26日

下午 11:34

$k(x_n, x_m)$  evaluated for all possible pairs  $x_n, x_m$

Maximum Margin Classifiers  $\rightarrow$  sparse

consider a two-class classification problem

$$y(x) = W^T \phi(x) + b \quad \phi(x): \text{input space} \rightarrow \text{feature space}$$

$x_1, \dots, x_N$  with labels  $t_1, \dots, t_N$   $t_n \in \{-1, 1\}$   
Assume linearly separable in feature space

$$y(x_n) > 0 \text{ for } t_n = +1$$

$$y(x_n) < 0 \text{ for } t_n = -1 \quad \Rightarrow \quad t_n y(x_n) > 0$$

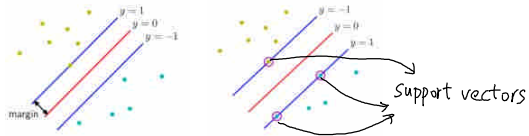
recall the perceptron algorithm, we might get multiple solutions.

we want to find the solution that give the smallest generalization error

SVMs achieve the goal by maximizing the margin

margin: the smallest distance between the decision boundary and any of the samples

$\Rightarrow$  choose the decision boundary of which the margin is maximum



2009年11月26日

下午 11:48

distance of a point  $x_n$  to the decision surface:

$$t_n \frac{W^T}{\|W\|} \left( \phi(x_n) - \phi(z) \right) \quad z \text{ is on the decision surface}$$

$$\Rightarrow W^T \phi(z) + b = 0$$

$$= \frac{t_n (W^T \phi(x_n) - W^T \phi(z))}{\|W\|}$$

$$= \frac{t_n (W^T \phi(x_n) + b)}{\|W\|} = \frac{t_n y(x_n)}{\|W\|}$$

maximize the margin

$$\arg \max_{w, b} \left\{ \frac{1}{\|W\|} \min_n [t_n (W^T \phi(x_n) + b)] \right\}$$

$$W \rightarrow kW \quad \text{scaling does not change } \frac{t_n y(x_n)}{\|W\|}$$

$$b \rightarrow kb$$

set  $t_n (W^T \phi(x^*) + b) = 1$  as a constraint  
 $x^*$  is the closest sample to the surface

so we have

$$t_n (W^T \phi(x_n) + b) \geq 1, \quad n = 1, \dots, N$$

equality holds  $\rightarrow$  active constraints

there must be at least two active constraints

we may rewrite the optimization problem as

$$\arg \max_{w, b} \frac{1}{\|W\|} \Rightarrow \arg \min_{w, b} \frac{1}{2} \|W\|^2 \quad \text{subject to}$$

$$t_n (W^T \phi(x_n) + b) \geq 1 \quad n = 1, \dots, N$$