

PRML学习笔记——第六章

- PRML学习笔记——第六章
 - Kernel Methods
 - 6.1 Dual Representations
 - 6.2 Constructing Kernels
 - 6.3. Radial Basis Function Networks
 - 6.3.1 Nadaraya-Watson model
 - 6.4. Gaussian Processes
 - 6.4.1 Gaussian processes for regression
 - 6.4.3 Learning the hyperparameters
 - 6.4.4 Automatic relevance determination
 - 6.4.5 Gaussian processes for classification
 - 6.4.6 Laplace approximation

Kernel Methods

6.1 Dual Representations

考虑SSE function:

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

求导令为0可得:

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

引入gram matrix $\mathbf{K} = \Phi \Phi^T$

那么SSE就能写成关于gram matrix的形式:

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}.$$

model的输出为:

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

可以看到,predict的时候只需要使用kernels的结果,避免使用原始data的feature $\phi(\mathbf{x})$.这在high,even infinity dimension的feature space上非常有用.

6.2 Constructing Kernels

一个valid kernel function等价于对应的gram matrix是半正定的.

利用这个等价关系仍然难以构造复杂的kernel function.

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

已知一些简单的kernel function,并利用这些性质,就可以构造复杂的kernel function.

6.3. Radial Basis Function Networks

将model内的basis function固定为radial basis function:

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

6.3.1 Nadaraya-Watson model

总的来说做的就是model $p(\mathbf{t}|\mathbf{x})$.

$\mathbb{E}[t|\mathbf{x}] = \sum_n k(\mathbf{x}, \mathbf{x}_n)t_n$,其中的 k 满足都大于0且求和是1.

6.4. Gaussian Processes

Gaussian Processes定义为a probability distribution over functions $y(\mathbf{x})$, such that the set of values of $y(\mathbf{x})$ evaluated at an arbitrary set of points x_1, \dots, x_N jointly have a Gaussian distribution.

6.4.1 Gaussian processes for regression

GPR可以看成是Bayes linear regression从weight space转到function space的extension,利用了kernel trick可以计算infinity dimension的basis function的情况.

Gaussian process定义的 \mathbf{y} :

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

其中的Covariance就是用Gram matrix给出的.利用第二章gaussian marginal的结果可以得到:

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

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

再次利用第二章中Gaussian joint的结果可以得到:

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

这就是我们需要的预测目标.重要的是Covariance可以用kernel trick直接给出结果.

6.4.3 Learning the hyperparameters

Gaussian process的predict部分依赖于covariance function.实际中常通过hyperparameter θ 来控制.一种方法确定 θ 是MLE给出point estimate:

$$\ln p(\mathbf{t} \mid \theta) = -\frac{1}{2}\ln|\mathbf{C}_N| - \frac{1}{2}\mathbf{t}^T \mathbf{C}_N^{-1} \mathbf{t} - \frac{N}{2}\ln(2\pi)$$

一般来说 $p(\mathbf{t} \mid \theta)$ 是nonconvex function,可以用一些基于gradient的optimize methods:

$$\frac{\partial}{\partial \theta_i} \ln p(\mathbf{t} \mid \theta) = -\frac{1}{2}\text{Tr}\left(\mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_i}\right) + \frac{1}{2}\mathbf{t}^T \mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_i} \mathbf{C}_N^{-1} \mathbf{t}.$$

另外也可以引入一个 θ 的prior来maximum posterior.还可以完全使用Bayes treatment,marginal over θ .

现在讨论的GPR model中的covariance内的 β 是constant,但在一些问题中可能dependent on \mathbf{x} .可以使用second gaussian process去model $\beta\mathbf{x}$

6.4.4 Automatic relevance determination

Target关于input variable(多个input的情况)可能具有不同程度的依赖强度.可以使用kernel function:

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp \left\{ -\frac{1}{2} \sum_{i=1}^D \eta_i (x_{ni} - x_{mi})^2 \right\} + \theta_2 + \theta_3 \sum_{i=1}^D x_{ni} x_{mi}$$

其中的 η 可以用来控制不同input variable对predict的重要性.这里面有Gaussian kernel,linear kernel和constant,通过 θ 可以实现soft *kernel selection*的效果.

6.4.5 Gaussian processes for classification

GP model给出的output是整个real axis,但对于classify problem,我们玩玩更需要一个在 $(0, 1)$ 区间的probability.

对于binary classify,我们可以在GP的输出后接一个logistic sigmoid function实现.记GP的输出为 a ,则target variable服从bernoulli distribution:

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

此时的predict function就变为了:

$$p(t_{N+1} = 1 | \mathbf{t}_N) = \int p(t_{N+1} = 1 | a_{N+1}) p(a_{N+1} | \mathbf{t}_N) da_{N+1}$$

其中 $p(t_{N+1} = 1 | a_{N+1}) = \sigma(a_{N+1})$.

note: 这个积分是intractable,所以需要用一些approximate approach.第四章用了sigmoid与Gaussian卷积近似,但由于GP的variable会随着input data增加而增加,并不能用中心极限定理近似成Gaussian.

6.4.6 Laplace approximation

对于上一节中的predict function,重点在确定积分中的 $p(a_{N+1} | \mathbf{t}_N)$

$$\begin{aligned} p(a_{N+1} | \mathbf{t}_N) &= \int p(a_{N+1}, \mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N \\ &= \frac{1}{p(\mathbf{t}_N)} \int p(a_{N+1}, \mathbf{a}_N) p(\mathbf{t}_N | a_{N+1}, \mathbf{a}_N) d\mathbf{a}_N \\ &= \frac{1}{p(\mathbf{t}_N)} \int p(a_{N+1} | \mathbf{a}_N) p(\mathbf{a}_N) p(\mathbf{t}_N | \mathbf{a}_N) d\mathbf{a}_N \\ &= \int p(a_{N+1} | \mathbf{a}_N) p(\mathbf{a}_N | \mathbf{t}_N) d\mathbf{a}_N \end{aligned}$$

而这里的 $p(a_{N+1} | \mathbf{a}_N)$ 在GPR中已经得到:

$$p(a_{N+1} | \mathbf{a}_N) = \mathcal{N}(a_{N+1} | \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{a}_N, c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k})$$

所以现在就需要用laplace approximate去近似 $p(\mathbf{t}_n | \mathbf{a}_n)$.

在近似完了之后就是一个对应的linear-Gaussian model.完了还有一部是要确定covariance function里的 θ .一个方法是用MLE,likelihood function为:

$$p(\mathbf{t}_N | \boldsymbol{\theta}) = \int p(\mathbf{t}_N | \mathbf{a}_N) p(\mathbf{a}_N | \boldsymbol{\theta}) d\mathbf{a}_N$$

这个积分仍然intractable,需要再一次laplace approximate.