

Series 3, 25 Oct 2019
(GPs; Model assessment and selection)

Teaching assistant: **Joanna Ficek**
joanna.ficek@inf.ethz.ch

Note: These are sample solutions. If you solved the problem in a different way it doesn't necessarily mean that your solution is wrong.

Solution 1 (Gaussian Processes cont.):

Solutions are based on Christopher M. Bishop, *Pattern Recognition and Machine Learning*. Springer Verlag (2006):

- a) 1. If $k_1(\mathbf{x}, \mathbf{x}')$ is a valid kernel, then there must exist a feature vector $\phi(\mathbf{x})$ such that

$$k_1(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') .$$

It follows that

$$ck_1(\mathbf{x}, \mathbf{x}') = \mathbf{u}(\mathbf{x})^T \mathbf{u}(\mathbf{x}')$$

where

$$\mathbf{u}(\mathbf{x}) = c^{1/2} \phi(\mathbf{x})$$

and so $ck_1(\mathbf{x}, \mathbf{x}')$ can be expressed as the scalar product of feature vectors, and hence is a valid kernel.

2. Similarly as above we can write

$$f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') = \mathbf{v}(\mathbf{x})^T \mathbf{v}(\mathbf{x}')$$

with

$$\mathbf{v}(\mathbf{x}) = f(\mathbf{x})\phi(\mathbf{x}) .$$

As before we can see that $f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$ can be expressed as the scalar product of feature vectors, and hence is a valid kernel.

3. We also know, that a necessary and sufficient condition for a function to be a valid kernel is that the Gram matrix \mathbf{K} , which elements are given by $k_1(\mathbf{x}, \mathbf{x}')$ should be positive semidefinite for all possible choices of the set $\{\mathbf{x}_n\}$. A matrix \mathbf{K} is positive semidefinite if, and only if,

$$\mathbf{a}^T \mathbf{K} \mathbf{a} \geq 0$$

for any choice of the vector \mathbf{a} . Let \mathbf{K}_1 be the Gram matrix for $k_1(\mathbf{x}, \mathbf{x}')$ and let \mathbf{K}_2 be the Gram matrix for $k_2(\mathbf{x}, \mathbf{x}')$. Then

$$\mathbf{a}^T (\mathbf{K}_1 + \mathbf{K}_2) \mathbf{a} = \mathbf{a}^T \mathbf{K}_1 \mathbf{a} + \mathbf{a}^T \mathbf{K}_2 \mathbf{a} \geq 0$$

where we have used the fact that \mathbf{K}_1 and \mathbf{K}_2 are positive semi-definite matrices, as well as the fact that the sum of two non-negative numbers will itself be non-negative. Thus, $k(\mathbf{x}, \mathbf{x}')$ defines a valid kernel.

4. Since we know that $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$ are valid kernels, we know that there exist mappings $\phi(\mathbf{x})$ and $\psi(\mathbf{x})$ such that

$$k_1(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

and

$$k_2(\mathbf{x}, \mathbf{x}') = \psi(\mathbf{x})^T \psi(\mathbf{x}').$$

Hence

$$\begin{aligned} k(\mathbf{x}, \mathbf{x}') &= k_1(\mathbf{x}, \mathbf{x}') k_2(\mathbf{x}, \mathbf{x}') \\ &= \phi(\mathbf{x})^T \phi(\mathbf{x}') \psi(\mathbf{x})^T \psi(\mathbf{x}') \\ &= \sum_{m=1}^M \phi_m(\mathbf{x}) \phi_m(\mathbf{x}') \sum_{n=1}^N \psi_n(\mathbf{x}) \psi_n(\mathbf{x}') \\ &= \sum_{m=1}^M \sum_{n=1}^N \phi_m(\mathbf{x}) \phi_m(\mathbf{x}') \psi_n(\mathbf{x}) \psi_n(\mathbf{x}') \\ &= \sum_{k=1}^K \psi_k(\mathbf{x}) \psi_k(\mathbf{x}') \\ &= \psi(\mathbf{x})^T \psi(\mathbf{x}') \end{aligned}$$

where $K = MN$ and

$$\psi_k(\mathbf{x}) = \phi_{(k-1) \oslash N + 1}(\mathbf{x}) \psi_{(k-1) \odot N + 1}(\mathbf{x})$$

where in turn \oslash and \odot denote integer division and remainder, respectively. Again we can see that $k(\mathbf{x}, \mathbf{x}')$ can be expressed as the scalar product of feature vectors, and hence is a valid kernel.

- b) The RBF kernel was used to generate the samples. σ corresponds to the vertical scaling and l to the horizontal scaling (for more information see the tutorial and <https://www.jgoertler.com/visual-exploration-gaussian-p>)
1. a:B ($\sigma = 0.8, l = 0.5$);
 2. c:A ($\sigma = 0.8, l = 2$);
 3. b:C ($\sigma = 0.33, l = 0.5$);

Solution 2 (Efficient Leave-One-Out Cross Validation):

The derivations include three steps:

Step 1.

$$\mathbf{w}_{(-i)}^* = \left(\mathbf{X}_{(-i)} \mathbf{X}_{(-i)}^T + \frac{(n-1)\mu}{2} \mathbf{I} \right)^{-1} \mathbf{X}_{(-i)} \mathbf{y}_{(-i)} \quad (1)$$

$$= (\mathbf{A} - \mathbf{x}_i \mathbf{x}_i^T)^{-1} (\mathbf{X} \mathbf{y} - \mathbf{x}_i y_i) \quad (2)$$

$$\stackrel{\text{SM}}{=} \mathbf{A}^{-1} (\mathbf{X} \mathbf{y} - \mathbf{x}_i y_i) + \frac{\mathbf{A}^{-1} \mathbf{x}_i \mathbf{x}_i^T \mathbf{A}^{-1} (\mathbf{X} \mathbf{y} - \mathbf{x}_i y_i)}{1 - \mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{x}_i} \quad (3)$$

$$= \mathbf{A}^{-1} \mathbf{X} \mathbf{y} - \mathbf{A}^{-1} \mathbf{x}_i \left(1 + \frac{\mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{x}_i}{1 - \mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{x}_i} \right) y_i + \frac{\mathbf{A}^{-1} \mathbf{x}_i \mathbf{x}_i^T \mathbf{A}^{-1} (\mathbf{X} \mathbf{y})}{1 - \mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{x}_i} \quad (4)$$

$$= \mathbf{A}^{-1} \mathbf{X} \mathbf{y} - \mathbf{A}^{-1} \mathbf{x}_i \left(\frac{1}{1 - \mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{x}_i} \right) y_i + \frac{\mathbf{A}^{-1} \mathbf{x}_i \mathbf{x}_i^T \mathbf{A}^{-1} (\mathbf{X} \mathbf{y})}{1 - \mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{x}_i}. \quad (5)$$

Step 2.

$$\mathbf{x}_i^T \mathbf{w}_{(-i)}^* \stackrel{\text{Eq. (5)}}{=} \mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{X} \mathbf{y} + \frac{\mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{x}_i \mathbf{x}_i^T \mathbf{A}^{-1} (\mathbf{X} \mathbf{y})}{1 - \mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{x}_i} - \left(\frac{1}{1 - \mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{x}_i} \right) \mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{x}_i y_i \quad (6)$$

$$= \left(1 + \frac{\mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{x}_i}{1 - \mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{x}_i} \right) \mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{X} \mathbf{y} - \left(\frac{1}{1 - \mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{x}_i} \right) \mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{x}_i y_i \quad (7)$$

$$= \left(\frac{1}{1 - \mathbf{x}_i^T \mathbf{A}^{-1} \mathbf{x}_i} \right) \mathbf{x}_i^T \mathbf{A}^{-1} (\mathbf{X} \mathbf{y} - \mathbf{x}_i y_i) \quad (8)$$

$$= \left(\frac{1}{1 - s_i} \right) (\hat{y}_i - s_i y_i). \quad (9)$$

Step 3.

$$y_i - \mathbf{x}_i^T \mathbf{w}_{(-i)}^* \stackrel{\text{Eq. (9)}}{=} y_i - \left(\frac{1}{1 - s_i} \right) (\hat{y}_i - s_i y_i) \quad (10)$$

$$= \left(1 + \frac{s_i}{1 - s_i} \right) y_i - \left(\frac{1}{1 - s_i} \right) \hat{y}_i \quad (11)$$

$$= \left(\frac{1}{1 - s_i} \right) (y_i - \hat{y}_i). \quad (12)$$

Solution 3 (Jackknife estimator):

a) The cumulative distribution function is

$$\mathbb{P}(X_{(n)} \leq x) = \mathbb{P}(X_1 \leq x, X_2 \leq x, \dots, X_n \leq x) = \prod_{i=1}^n \mathbb{P}(X_i \leq x) = \left(\frac{x}{\theta} \right)^n.$$

Hence we can compute the probability density function (PDF) as

$$p_{(n)}(x) = \frac{d}{dx} \mathbb{P}(X_{(n)} \leq x) = \frac{n}{\theta} \left(\frac{x}{\theta} \right)^{n-1}.$$

The average of the estimator is then

$$\mathbb{E}[X_{(n)}] = \int_0^\theta dx \, x \frac{n}{\theta} \left(\frac{x}{\theta} \right)^{n-1} = \frac{n}{n+1} \theta$$

One can see that the estimator $X_{(n)}$ underestimates the bound θ .

b) The replicate estimator $\hat{S}_{n-1}^{(-i)}$ is the maximum of the samples $\{X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n\}$. Let assume $i^* = \arg \max_i X_i$. Then the replicate estimator reads

$$\hat{S}_{n-1}^{(-i)} = \begin{cases} X_{(n)} & i \neq i^* \\ X_{(n-1)} & i = i^* \end{cases}$$

c) The Jackknife estimator reads

$$\begin{aligned} \hat{S}_n^{JK} &= \hat{S}_n - (n-1) \left(\frac{1}{n} \sum_{i=1}^n \hat{S}_{n-1}^{(-i)} - \hat{S}_n \right) = X_{(n)} - (n-1) \left(\frac{n-1}{n} X_{(n)} + \frac{1}{n} X_{(n-1)} - X_{(n)} \right) \\ &= X_{(n)} + \frac{n-1}{n} (X_{(n)} - X_{(n-1)}). \end{aligned}$$

Hence, the Jackknife estimator modifies the estimator \hat{S}_n by adding a positive correction.

d) The cumulative distribution function (CDF) of $X_{(n)}$ is

$$\begin{aligned}\mathbb{P}(X_{(n-1)} \leq x) &= \mathbb{P}(X_{(n-1)} \leq x, X_{(n)} \leq x) + \mathbb{P}(X_{(n-1)} \leq x, X_{(n)} \geq x) \\ &= \mathbb{P}(X_{(n)} \leq x) + \sum_{i=1}^n \mathbb{P}(X_i \geq x, X_1 \leq x, \dots, X_{i-1} \leq x, X_{i+1} \leq x, \dots, X_n \leq x) \\ &= \left(\frac{x}{\theta}\right)^n + n \frac{\theta - x}{\theta} \left(\frac{x}{\theta}\right)^{n-1}.\end{aligned}$$

Therefore, the probability density function (PDF) is

$$p_{(n-1)}(x) = \frac{d}{dx} \mathbb{P}(X_{(n-1)} \leq x) = \frac{n(n-1)}{\theta} \left(\frac{x}{\theta}\right)^{n-2} \left(1 - \frac{x}{\theta}\right).$$

We can now compute the expected value as

$$\mathbb{E}[X_{(n-1)}] = \int_0^\theta x p_{(n-1)}(x) dx = \int_0^\theta x \frac{n(n-1)}{\theta} \left(\frac{x}{\theta}\right)^{n-2} \left(1 - \frac{x}{\theta}\right) dx = \frac{n-1}{n+1} \theta,$$

and, finally,

$$\mathbb{E}[\hat{S}_n^{JK}] = \left(1 - \frac{1}{n^2 + n}\right) \theta$$

Note that the bias of the Jackknife estimator \hat{S}_n^{JK} is smaller than the bias of the original estimator \hat{S}_n by a factor of n . Also note that the Jackknife estimator does not exploit any property of the distribution, and thus it can be used also in cases where the actual distribution is unknown.

Solution 4 (Model selection: Bayesian Information Criterion):

BIC can be seen as a large n approximation to the log model evidence and hence, in the below solution we consider $n \rightarrow \infty$. We use the Laplace approximation to the log model evidence around the mode of the posterior distribution $p(\theta^{(m)} | \mathcal{D}^{(n)})$

$$\ln p(\mathcal{D}^{(n)}) \approx \ln p(\mathcal{D}^{(n)} | \theta_{\text{MAP}}^{(m)}) + \ln p(\theta_{\text{MAP}}^{(m)}) + \frac{m}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{A}|,$$

where $\mathbf{A} = -\frac{\partial^2}{\partial \theta^{(m)} \partial \theta^{(m)}} \ln p(\theta_{\text{MAP}}^{(m)} | \mathcal{D}^{(n)})$ is the Hessian of the minus log posterior at $\theta_{\text{MAP}}^{(m)}$.

We neglect the log prior $\ln p(\theta_{\text{MAP}}^{(m)})$ at $\theta_{\text{MAP}}^{(m)}$ and consider a simple case where the Hessian $\mathbf{A} \in \mathbb{R}^{m \times m}$ is a diagonal matrix. We assume that the Hessian has a full rank. In addition, we assume the data to be iid, which allows us to write the likelihood as a product $p(\mathcal{D}^{(n)} | \theta_{\text{MAP}}^{(m)}) = \prod_{i=1}^n p(\mathbf{x}_i | \theta_{\text{MAP}}^{(m)})$, and therefore,

$$[\mathbf{A}]_{jj} = \frac{\partial^2}{\partial \theta_j^2} \ln p(\theta_{\text{MAP}}^{(m)} | \mathcal{D}^{(n)}) \quad (13)$$

$$= \frac{\partial^2}{\partial \theta_j^2} \left[\ln p(\mathcal{D}^{(n)} | \theta_{\text{MAP}}^{(m)}) + \ln p(\theta_{\text{MAP}}^{(m)}) - \ln p(\mathcal{D}^{(n)}) \right] \quad (14)$$

$$= \frac{\partial^2}{\partial \theta_j^2} \left[\ln p(\mathcal{D}^{(n)} | \theta_{\text{MAP}}^{(m)}) + \ln p(\theta_{\text{MAP}}^{(m)}) \right] \quad (15)$$

$$= \sum_{i=1}^n \frac{\partial^2}{\partial \theta_j^2} \ln p(\mathbf{x}_i | \theta_{\text{MAP}}^{(m)}) + \frac{\partial^2}{\partial \theta_j^2} \ln p(\theta_{\text{MAP}}^{(m)}) \quad (16)$$

$$\sim nc_j. \quad (17)$$

Thus, $|\mathbf{A}| \sim n^m \prod_j c_j$, which leads us to the following result

$$\ln |\mathbf{A}| \sim m \ln n.$$

Finally, one can write

$$-2 \ln p(\mathcal{D}^{(n)}) \approx -2 \ln p(\mathcal{D}^{(n)} | \theta^{(m)}) + m \ln n = \text{BIC}(\mathcal{D}^{(n)}).$$

Another way of solving the exercise can be found in Kevin P. Murphy, *Machine Learning: A Probabilistic Perspective*. MIT Press (2012), p.255, 256.