

# MML lecture extra notes, Nov 26, 2021

The contents in this note is non-examinable, but you might find it useful as preparation for the “Probabilistic Inference” course in Spring.

## Full derivations for kernel PCA

In non-linear PCA, we wish to find orthonormal projection directions  $\{\mathbf{b}_1, \dots, \mathbf{b}_M\}$  for the non-linear features  $\{\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N)\}$ ,  $\phi(\mathbf{x}) \in \mathbb{R}^{p \times 1}$  such that  $\sum_{m=1}^M \mathbb{V}[\phi(\mathbf{x}_n)^\top \mathbf{b}_m]$  is maximised. Here we assume that the dataset  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  is pre-processed such that  $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\phi(\mathbf{x})] = \mathbf{0}$ . From previous discussions of PCA we see that the optimal solution of  $\{\mathbf{b}_1, \dots, \mathbf{b}_m\}$  is the leading  $M$  eigenvectors of matrix  $\mathbf{S} = \frac{1}{N} \sum_{n=1}^n \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \in \mathbb{R}^{p \times p}$ .

Now consider eigendecomposition of  $\mathbf{S}$ , which requires solving for  $\mathbf{b}_m$  and  $\lambda_m$  the following equation:

$$\frac{1}{N} \sum_{n=1}^n \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \mathbf{b}_m = \lambda_m \mathbf{b}_m. \quad (1)$$

Notice that both  $\lambda_m$  and  $\phi(\mathbf{x}_n)^\top \mathbf{b}_m$  are scalars. So the optimal solution for  $\mathbf{b}_m$  has the following form

$$\mathbf{b}_m := \sum_{n=1}^N \alpha_{mn} \phi(\mathbf{x}_n) \quad (2)$$

for a particular set of  $\{\alpha_{mn}\}$  parameters. In the eigendecomposition problem (1) it means  $\alpha_{mn} = \frac{\phi(\mathbf{x}_n)^\top \mathbf{b}_m}{N \lambda_m}$ . But instead of solving for  $\lambda_m$  directly (as an eigenvalue of  $\mathbf{S}$ ) which typically has  $\mathcal{O}(p^3)$  computational cost, we can directly use the alternative parameterisation (2) and compute the objective for  $\mathbf{b}_m$  as follows:

$$\begin{aligned} \mathbf{b}_m^\top \mathbf{S} \mathbf{b}_m &= \frac{1}{N} \sum_{n=1}^n \mathbf{b}_m^\top \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \mathbf{b}_m \\ &= \frac{1}{N} \sum_{n=1}^n \left( \phi(\mathbf{x}_n)^\top \sum_{n'=1}^N \alpha_{mn'} \phi(\mathbf{x}_{n'}) \right)^2 \\ &= \frac{1}{N} \sum_{n=1}^n \left( \sum_{n'=1}^N \alpha_{mn'} \mathcal{K}(\mathbf{x}_n, \mathbf{x}_{n'}) \right)^2. \end{aligned} \quad (3)$$

Now define  $\boldsymbol{\alpha}_m = [\alpha_{m1}, \dots, \alpha_{mN}]^\top$ , and define the kernel matrix as  $\mathbf{K}_{\mathbf{X}\mathbf{X}} \in \mathbb{R}^{N \times N}$  with entries  $\mathbf{K}_{\mathbf{X}\mathbf{X}}[i, j] = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$ . This allows us to further rewrite the objective as

$$\mathbf{b}_m^\top \mathbf{S} \mathbf{b}_m = \frac{1}{N} \sum_{n=1}^n (\mathbf{K}_{\mathbf{X}\mathbf{X}}[n, \cdot] \boldsymbol{\alpha}_m)^2 = \frac{1}{N} \|\mathbf{K}_{\mathbf{X}\mathbf{X}} \boldsymbol{\alpha}_m\|_2^2. \quad (4)$$

For the constraints on  $\mathbf{b}_m$ , they also translates to the constraints on  $\boldsymbol{\alpha}_m$ :

$$\begin{aligned} \|\mathbf{b}_m\|_2^2 = 1 &\Rightarrow \sum_{n, n'} \alpha_{mn} \alpha_{mn'} \mathcal{K}(\mathbf{x}_n, \mathbf{x}_{n'}) = \boldsymbol{\alpha}_m^\top \mathbf{K}_{\mathbf{X}\mathbf{X}} \boldsymbol{\alpha}_m = 1, \\ \mathbf{b}_i \perp \mathbf{b}_j, \forall i \neq j &\Rightarrow \sum_{n, n'} \alpha_{in} \alpha_{jn'} \mathcal{K}(\mathbf{x}_n, \mathbf{x}_{n'}) = \boldsymbol{\alpha}_i^\top \mathbf{K}_{\mathbf{X}\mathbf{X}} \boldsymbol{\alpha}_j = 0. \end{aligned} \quad (5)$$

Readers are encouraged to solve this problem themselves using the Lagrange multiplier method. Such method returns solutions  $\{\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_M\}$  as the leading  $M$  eigenvectors of matrix  $\frac{1}{N} \mathbf{K}_{\mathbf{X}\mathbf{X}} \in \mathbb{R}^{N \times N}$ . In prediction, given a new datapoint  $\hat{\mathbf{x}}$ , dimension reduction can be done by

$$\mathbf{b}_m^\top \phi(\hat{\mathbf{x}}) = \sum_{n=1}^N \alpha_{mn} \mathcal{K}(\hat{\mathbf{x}}, \mathbf{x}_n) = \mathbf{K}_{\hat{\mathbf{x}}\mathbf{X}} \boldsymbol{\alpha}_m. \quad (6)$$

Compare with performing eigendecomposition for  $\mathbf{S}$  directly, eigendecomposition of  $\mathbf{K}_{\mathbf{X}\mathbf{X}}$  typically has cost  $\mathcal{O}(N^3)$ , resulting in significant speed-up if  $N \ll p$  (or even when  $p \Rightarrow +\infty$ ). The kernel PCA can also be extended to settings where  $\phi(\mathbf{x})$  lies in a non-Euclidean space, see the next section for further discussions.

## From Bayesian linear regression to Gaussian processes

Recall the parametric non-linear regression model:

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y; \phi(\mathbf{x})^\top \boldsymbol{\theta}, \sigma^2), \quad p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}; \mathbf{0}, \lambda^{-1} \mathbf{I})$$

Now we derive a Gaussian process model from the above Bayesian linear regression model by lifting the prior from parameter space to function space. Recall that  $f(\mathbf{x}) = \phi(\mathbf{x})^\top \boldsymbol{\theta}$ . This means for any  $N > 0$  and  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times d}$ , we have

$$\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^\top = \Phi \boldsymbol{\theta} \in \mathbb{R}^{N \times 1}, \quad \Phi = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N)]^\top.$$

Also notice that for a Gaussian variable  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \Sigma)$ , we have  $\mathbf{A}\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{A}\Sigma\mathbf{A}^\top)$ . This means we can lift the prior on  $\boldsymbol{\theta}$  to a prior on  $\mathbf{f}$ :

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}; \mathbf{0}, \lambda^{-1} \mathbf{I}) \Rightarrow p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{f}; \mathbf{0}, \mathbf{K}_{\mathbf{X}\mathbf{X}}), \quad \mathbf{K}_{\mathbf{X}\mathbf{X}} = \lambda^{-1} \Phi \Phi^\top.$$

Importantly, this lifting procedure  $\mathbf{f} = \Phi \boldsymbol{\theta}$  makes  $f(\mathbf{x}_i)$  and  $f(\mathbf{x}_j)$  correlated. Since we can perform such lifting for any  $N > 0$ , this results in a Gaussian process (GP) prior written as follows:

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}; \mathbf{0}, \lambda^{-1} \mathbf{I}) \Rightarrow f(\cdot) \sim \mathcal{GP}(0, \mathcal{K}(\cdot, \cdot)), \quad \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \lambda^{-1} \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j).$$

For the likelihood terms we can perform similar lifting as well:

$$\prod_{n=1}^N p(y_n|\mathbf{x}_n, \boldsymbol{\theta}) = \prod_{n=1}^N \mathcal{N}(y_n; \phi(\mathbf{x}_n)^\top \boldsymbol{\theta}, \sigma^2) = \mathcal{N}(\mathbf{y}; \Phi \boldsymbol{\theta}, \sigma^2 \mathbf{I}) \Rightarrow p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}; \mathbf{f}, \sigma^2 \mathbf{I}).$$

Combining both results returns a GP regression model as

$$p(y|f, \mathbf{x}) = \mathcal{N}(y; f(\mathbf{x}), \sigma^2), \quad f(\cdot) \sim \mathcal{GP}(0, \mathcal{K}(\cdot, \cdot)).$$

In general  $\phi(\mathbf{x}_i)$  might be infinite-dimensional (or even in non-Euclidean space), so one would rather directly define the kernel function  $\mathcal{K}(\cdot, \cdot)$  instead of specifying the form of  $\phi(\cdot)$ . Specifically, if  $\mathcal{K}(\cdot, \cdot)$  is a *reproducing kernel* of a *reproducing kernel Hilbert space* (RKHS)  $\mathcal{H}$  equipped with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ , then one can use  $\phi(\mathbf{x}) = \mathcal{K}(\mathbf{x}, \cdot)$  as the non-linear feature (which is a function so is infinite-dimensional) and compute the inner product as  $\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle_{\mathcal{H}} = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$ .