

Manifold Learning for Signal and Visual Processing

Lecture 10: Introduction to Gaussian Processes

Radu Horaud
INRIA Grenoble Rhone-Alpes, France
Radu.Horaud@inria.fr
<https://team.inria.fr/perception>

Outline of This Lecture

- Back to linear regression.
- Bayesian linear regression.
- Gaussian process problem statement.
- Gaussian process for regression.

Material for This Lecture

- C. M. Bishop (2006). Pattern Recognition and Machine Learning. Chapters 3 and 6 (section 6.4).
- More involved readings:
 - Rasmussen and Williams (2006). Gaussian Processes for Machine Learning.

Linear Regression

- The task is to estimate the parameters $\mathbf{w} = (w_0, \dots, w_{M-1})^\top$ in:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x})$$

- from a training dataset of N observations $\{\mathbf{x}_n\}$ together with corresponding target values $\{t_n\}$.
- The target value is given by $t = y(\mathbf{x}, \mathbf{w}) + \epsilon$ where ϵ is a zero-mean Gaussian variable with precision β . Thus we can write:

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

Maximum Likelihood Estimation

- The log-likelihood function is, with $\mathbf{t} = (t_1, \dots, t_n)$:

$$\ln p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1})$$

- MLE solution:

$$\mathbf{w}_{\text{ML}} = (\boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{t}$$

- and the variance (inverse of precision):

$$\beta_{\text{ML}}^{-1} = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{w}_{\text{ML}}^\top \boldsymbol{\phi}(\mathbf{x}_n))^2$$

Bayesian Linear Regression

- Prior distribution of the model parameters:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$$

- The posterior distribution is proportional to the product of the likelihood function $p(\mathbf{t}|\mathbf{w})$ and the prior $p(\mathbf{w})$. This results in (see chapter 2 of Bishop, eq. (2.116)):

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$$

- with:

$$\begin{aligned}\mathbf{m}_N &= \beta \mathbf{S}_N \Phi^T \mathbf{t} \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \Phi^T \Phi.\end{aligned}$$

Maximization of the Posterior Distribution

$$\ln p(\mathbf{t}|\mathbf{w}) = -\frac{\beta}{2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \phi(\mathbf{x}_n))^2 - \frac{\alpha}{2} \mathbf{w}^\top \mathbf{w} + \text{const.}$$

- This is equivalent to a regularized least-square optimization problem with respect to \mathbf{w} and with a regularization coefficient $\lambda = \alpha/\beta$ (see Bishop section 3.14).

Predictive Distribution

- Predict a value for t for a new value of \mathbf{x} given the model parameters \mathbf{w} (I simplified the Bishop notations):

$$p(t) \propto \int p(t|\mathbf{w})p(\mathbf{w})d\mathbf{w}$$

- Using the previous formulae we obtain the *predictive distribution* $p(t|\mathbf{w})$:

$$p(t|\mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t, \mathbf{m}_N^\top \phi(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$
$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \phi(\mathbf{x})^\top \mathbf{S}_N \phi(\mathbf{x}).$$

Linear Regression Again

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

- Let $\mathbf{y} = \{y_n\}$ with $y_n = y(\mathbf{x}_n)$ and let Φ be the following $N \times M$ matrix:

$$\Phi = \begin{bmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \vdots & \vdots & & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \dots & \phi_{M-1}(\mathbf{x}_N) \end{bmatrix}$$

- Altogether this can be written:

$$\mathbf{y} = \Phi \mathbf{w}$$

Probability Distribution of \mathbf{y}

- We seek the joint probability distribution of $\mathbf{y} = \{y_1, \dots, y_N\}$ using the prior $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$.
- \mathbf{y} is a linear combination of Gaussian variables, it is Gaussian itself, with:

$$E[\mathbf{y}] = \Phi E[\mathbf{w}] = \mathbf{0}$$

$$\text{cov}[\mathbf{y}] = E[\mathbf{y}\mathbf{y}^\top] = \Phi E[\mathbf{w}\mathbf{w}^\top] \Phi^\top = \frac{1}{\alpha} \Phi \Phi^\top = \mathbf{K}$$

- The matrix \mathbf{K} is a Gramm (or kernel) matrix with elements given by:

$$K_{nm} = \frac{1}{\alpha} \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$$

Gaussian Process Definition (Bishop)

- A Gaussian process (GP) is defined as probability distribution over functions $y(\mathbf{x})$ such that the set of values y_1, \dots, y_N evaluated at an arbitrary set of points $\mathbf{x}_1, \dots, \mathbf{x}_N$, jointly have a Gaussian distribution.
- When $\mathbf{x} \in \mathbb{R}^2$ this is known as a Gaussian random field.
- A stochastic process $y(\mathbf{x})$ is specified by the joint probability distribution for a finite set of values in a consistent manner.
- Key point about stochastic GP: the joint distribution $p(y_1, \dots, y_N)$ is specified completely by the second-order statistics (mean and covariance)

GP for Regression

- Let's consider again that the observed target values have Gaussian noise, namely $t_n = y_n + \epsilon_n$, or:

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

- The joint distribution of the target values:

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

- From above, we can write the distribution of \mathbf{y} :

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

where \mathbf{K} is a kernel matrix. The kernel must be chosen such as to express the fact that if \mathbf{x}_n and \mathbf{x}_m are similar the corresponding values of $y(\mathbf{x}_n)$ and $y(\mathbf{x}_m)$ will be strongly correlated.

GP for Regression

- The marginal distribution $p(\mathbf{t})$, conditioned on the observed input values $\{\mathbf{x}_n\}$, can be found by integrating over \mathbf{y} and using the results of chapter 2 (Bishop):

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

- with:

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

A Gaussian Process Example

- A “widely” used kernel function:

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left(-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2\right) + \theta_2 + \theta_3 \mathbf{x}_n^\top \mathbf{x}_m$$

- $\boldsymbol{\theta} = (\theta_0, \theta_1, \theta_2, \theta_3)$ are the hyper-parameters of the GP model.
- We found an expression for $p(\mathbf{t})$ but it would be more correct to write $p(\mathbf{t}|\boldsymbol{\theta})$.

The Predictive Distribution

- Let \mathbf{x}_{N+1} be a new input vector and we want to predict the target t_{N+1} . We need to compute the predictive distribution $p(t_{N+1}|\mathbf{t}_N, \mathbf{X}_{N+1})$ with $\mathbf{X}_{N+1} = (\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{x}_{N+1})$, for convenience we drop \mathbf{X} .
- Hence the predictive distribution is written $p(t_{N+1}|\mathbf{t}_N)$. It can be obtained from the joint distribution:

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

- The covariance matrix can be decomposed such as to use the covariance that was estimated from the training data, namely:

$$\mathbf{C}_{N+1} = \begin{bmatrix} & & & k(\mathbf{x}_1, \mathbf{x}_{N+1}) \\ & & & \vdots \\ & & & \vdots \\ & & & k(\mathbf{x}_N, \mathbf{x}_{N+1}) \\ & & & \underbrace{k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1}}_c \\ & & \underbrace{k(\mathbf{x}_1, \mathbf{x}_{N+1}) \quad \dots \quad k(\mathbf{x}_N, \mathbf{x}_{N+1})}_{\mathbf{k}^\top} & \end{bmatrix}$$

The Predictive Distribution

- Using various results (Bishop section 2, eqs. (2.81) and (2.82)) we see that the predictive (conditional) distribution is Gaussian with mean and covariance:

$$p(t_{N+1}|\mathbf{t}_N) = \mathcal{N}(t_{N+1} | \underbrace{\mathbf{k}^\top \mathbf{C}_N^{-1} \mathbf{t}_N}_{\text{mean}}, \underbrace{c - \mathbf{k}^\top \mathbf{C}_N^{-1} \mathbf{k}}_{\text{cov}})$$

Remarks

- The training involves the estimation of the hyper-parameters θ by maximizing the log-likelihood function $\ln p(\mathbf{t}_N|\theta)$ which is a non-convex problem.
- It is possible to introduce a prior $p(\theta)$ and maximize the log-posterior using gradient-based methods. This implies marginalization, which is intractable in the general case.
- The covariance matrix \mathbf{C}_N of size $N \times N$ must be inverted.
- There is no general methodology to select the kernel function!