

Gaussian Process Regression

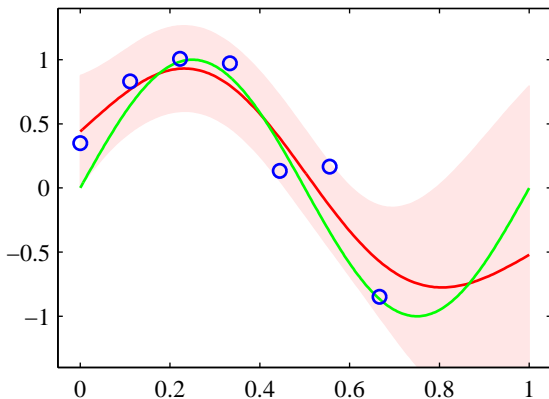
Seungjin Choi

Department of Computer Science and Engineering
Pohang University of Science and Technology
77 Cheongam-ro, Nam-gu, Pohang 37673, Korea
seungjin@postech.ac.kr

Pictorial Illustration of GP Regression

Green curve: the true sinusoidal function from which the data points are obtained by sampling and addition of Gaussian noise.

Red line: the mean of GP predictive distribution and the shaded region corresponds to plus and minus two standard deviations.



Gaussian Processes

Definition: A **Gaussian Process (GP)** is a collection of random variables, any finite number of which has a joint Gaussian distribution.

- ▶ A Gaussian process is a generalization of a multivariate Gaussian distribution to **infinitely many variables**.
- ▶ GP defines a **a distribution over functions** of the form $f : \mathcal{X} \mapsto \mathbb{R}$, which is completely specified by **mean function $\mu(\mathbf{x})$** and **covariance function $k(\mathbf{x}, \mathbf{x}')$** :

$$f(\mathbf{x}) \sim \mathcal{GP}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')),$$

where

$$\begin{aligned}\mu(\mathbf{x}) &= \mathbb{E}\{f(\mathbf{x})\}, \\ k(\mathbf{x}, \mathbf{x}') &= \mathbb{E}\{(f(\mathbf{x}) - \mu(\mathbf{x})) (f(\mathbf{x}') - \mu(\mathbf{x}'))\} \\ &= l_f \exp\left\{-\frac{1}{2\rho}(\mathbf{x} - \mathbf{x}')^\top \mathbf{L}(\mathbf{x} - \mathbf{x}')\right\},\end{aligned}$$

which is referred to as **squared exponential kernel**, $\mathbf{L} = \text{diag}(\mathbf{I})$, $[I]_i$ is a hyperparameter to determine a relevance of the i th input dimension.

Regression

- ▶ The **regression model** is described by

$$y_t = f(\mathbf{x}_t) + \epsilon_t.$$

- ▶ **Data:** $\mathcal{D} = \{(\mathbf{X}, \mathbf{y})\}$

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N},$$

$$\mathbf{y} = [y_1, \dots, y_N]^T \in \mathbb{R}^N.$$

- ▶ **Goal:** Estimate the regression function $f(\cdot)$ to make prediction of $y_* = f(\mathbf{x}_*)$.
 - ▶ **Parametric:** Maximum likelihood, Bayesian inference (**parametric model** $f_w(\mathbf{x})$).
 - ▶ **Nonparametric:** Kernel regression (Nadaraya-Watson estimator), GP regression (**let the data explain, $f(\mathbf{x})$**).

Maximum Likelihood (Parametric Model)

- ▶ Model:

$$y_t = f_w(\mathbf{x}_t) + \epsilon_t.$$

- ▶ Gaussian likelihood:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) \propto \prod_{t=1}^N \exp \left\{ -\frac{1}{2\sigma^2} (y_t - f_w(\mathbf{x}_t))^2 \right\}.$$

- ▶ Maximize the log-likelihood:

$$\mathbf{w}_{ML} = \arg \max_w \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}).$$

- ▶ Make predictions by plugging in the ML estimate:

$$p(y_*|\mathbf{x}_*, \mathbf{w}_{ML}).$$

Bayesian Inference (Parametric Model)

- ▶ Model:

$$y_t = f_w(\mathbf{x}_t) + \epsilon_t.$$

- ▶ Gaussian likelihood:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) \propto \prod_{t=1}^N \exp \left\{ -\frac{1}{2\sigma^2} (y_t - f_w(\mathbf{x}_t))^2 \right\}.$$

- ▶ Prior over parameters:

$$p(\mathbf{w}|\lambda) \quad (\text{hyperparameters } \lambda).$$

- ▶ Posterior over parameters:

$$p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \lambda) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w}|\lambda)}{p(\mathbf{y}|\mathbf{X})}.$$

Bayesian Inference (Parametric Model)

- ▶ Make predictions by computing predictive distribution:

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}, \lambda) = \int p(y_* | \mathbf{w}, \mathbf{x}_*) p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \lambda) d\mathbf{w}.$$

- ▶ Marginal likelihood:

$$p(\mathbf{y} | \mathbf{X}) = \int p(\mathbf{y} | \mathbf{X}, \mathbf{w}) p(\mathbf{w} | \lambda) d\mathbf{w}.$$

- ▶ Estimate hyperparameters by maximizing the log marginal likelihood:

$$\hat{\lambda} = \arg \max_{\lambda} \log p(\mathbf{y} | \mathbf{X}).$$

Kernel Regression

- ▶ Model:

$$y_t = f(\mathbf{x}_t) + \epsilon_t.$$

- ▶ Nadayara-Watson estimator:

$$f(\mathbf{x}_*) = \frac{\sum_{t=1}^N y_t k(\mathbf{x}_*, \mathbf{x}_t, \lambda)}{\sum_{l=1}^N k(\mathbf{x}_*, \mathbf{x}_l, \lambda)},$$

where

$$k(\mathbf{x}, \mathbf{x}_t, \lambda) = \frac{1}{Z} \exp \{-\lambda \|\mathbf{x} - \mathbf{x}_t\|^2\}.$$

- ▶ Computes the weighted average of y_t 's near \mathbf{x}_* .

Nadaya-Watson Estimator (Detailed Derivation)

Recall

$$\begin{aligned}f(\mathbf{x}) &= \mathbb{E}\{y|\mathbf{x}\} \\&= \int y p(y|\mathbf{x}) dy \\&= \int y \frac{p(\mathbf{x}, y)}{p(\mathbf{x})} dy.\end{aligned}$$

Use [kernel density estimation](#) to determine both $p(\mathbf{x}, y)$ and $p(\mathbf{x})$:

$$\begin{aligned}\hat{p}(\mathbf{x}, y) &= \frac{1}{N} \sum_{t=1}^N k(\mathbf{x}, \mathbf{x}_t, \lambda_x) k(y, y_t, \lambda_y), \\ \hat{p}(\mathbf{x}) &= \frac{1}{N} \sum_{t=1}^N k(\mathbf{x}, \mathbf{x}_t, \lambda_x).\end{aligned}$$

Compute

$$\begin{aligned}\int y \hat{p}(\mathbf{x}, y) dy &= \frac{1}{N} \sum_{t=1}^N \int k(\mathbf{x}, \mathbf{x}_t, \lambda_x) y k(y, y_t, \lambda_y) dy \\ &= \frac{1}{N} \sum_{t=1}^N k(\mathbf{x}, \mathbf{x}_t, \lambda_x) \underbrace{\int y \frac{1}{Z_y} \exp\{-\lambda_y(y - y_t)^2\} dy}_{y_t} \\ &= \frac{1}{N} \sum_{t=1}^N y_t k(\mathbf{x}, \mathbf{x}_t, \lambda_x)\end{aligned}$$

Therefore,

$$\begin{aligned}f(\mathbf{x}) &= \frac{\int y \hat{p}(\mathbf{x}, y) dy}{\hat{p}(\mathbf{x})} \\ &= \frac{\sum_{t=1}^N y_t k(\mathbf{x}, \mathbf{x}_t, \lambda_x)}{\sum_{l=1}^N k(\mathbf{x}, \mathbf{x}_l, \lambda_x)}.\end{aligned}$$

Gaussian Process Regression

- ▶ Model:

$$y_t = f(\mathbf{x}_t) + \epsilon_t,$$

where $f(\cdot)$ is referred to as **latent function**.

- ▶ Latent vector:

$$\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^\top \in \mathbb{R}^N.$$

Note that "parameters" are function itself in GPR model.

- ▶ Gaussian likelihood:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I}).$$

- ▶ Gaussian process prior (zero mean):

$$\begin{aligned} f(\mathbf{x}) &\sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}')), \\ p(\mathbf{f}|\mathbf{X}) &= \mathcal{N}(\mathbf{f}|0, \mathbf{K}). \end{aligned}$$

GP Regression (Cont'd)

- ▶ Gaussian process posterior:

$$f(\mathbf{x})|\mathbf{X}, \mathbf{y} \sim \mathcal{GP}(\bar{\mu}(\mathbf{x}), \bar{k}(\mathbf{x}, \mathbf{x}')),$$

where

$$\bar{\mu}(\mathbf{x}) = k(\mathbf{x}, \mathbf{X}) \left[k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I} \right]^{-1} \mathbf{y},$$

$$\bar{k}(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}, \mathbf{x}') - k(\mathbf{x}, \mathbf{X}) \left[k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I} \right]^{-1} k(\mathbf{X}, \mathbf{x}).$$

- ▶ Predictive distribution:

$$\begin{aligned} p(y_*|\mathbf{X}, \mathbf{y}, \mathbf{x}_*) &= \mathcal{N} \left(y_* \mid k(\mathbf{x}_*, \mathbf{X}) \left[k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I} \right]^{-1} \mathbf{y}, \right. \\ &\quad \left. k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X}) \left[k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I} \right]^{-1} k(\mathbf{X}, \mathbf{x}_*) + \sigma^2 \right), \end{aligned}$$

where

$$k(\mathbf{x}_*, \mathbf{X}) = [k(\mathbf{x}_*, \mathbf{x}_1), \dots, k(\mathbf{x}_*, \mathbf{x}_N)] \in \mathbb{R}^{1 \times N},$$

$$k(\mathbf{X}, \mathbf{X}) = [k(\mathbf{x}_i, \mathbf{x}_j)] \in \mathbb{R}^{N \times N}.$$

GP Regression: Detailed Derivation

Let $\mathbf{f}_* \in \mathbb{R}^T$ be latent function values evaluated at test data points
 $\mathbf{X}_* \in \mathbb{R}^{D \times T}$.

We first write the joint distribution of the observed target values and the function values at the test locations under the prior:

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I} & k(\mathbf{X}, \mathbf{X}_*) \\ k(\mathbf{X}_*, \mathbf{X}) & k(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right),$$

It follows from the Gaussian Identity that we have

$$\begin{aligned} \mathbf{f}_* | \mathbf{y}, \mathbf{X}, \mathbf{X}_* &\sim \mathcal{N} \left(\mathbf{f}_* | k(\mathbf{X}_*, \mathbf{X}) \left[k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N \right]^{-1} \mathbf{y}, \right. \\ &\quad \left. k(\mathbf{X}_*, \mathbf{X}_*) - k(\mathbf{X}_*, \mathbf{X}) \left[k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N \right]^{-1} k(\mathbf{X}, \mathbf{X}_*) \right), \end{aligned}$$

leading to

$$\begin{aligned} p(\mathbf{y}_* | \mathbf{y}, \mathbf{X}, \mathbf{X}_*) &= \mathcal{N} \left(\mathbf{f}_* | k(\mathbf{X}_*, \mathbf{X}) \left[k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N \right]^{-1} \mathbf{y}, \right. \\ &\quad \left. k(\mathbf{X}_*, \mathbf{X}_*) - k(\mathbf{X}_*, \mathbf{X}) \left[k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N \right]^{-1} k(\mathbf{X}, \mathbf{X}_*) + \sigma^2 \mathbf{I}_T \right). \end{aligned}$$

References

- ▶ C. E. Rasmussen and C. K. I. Williams, "Gaussian Processes for Machine Learning," MIT Press, 2006.
- ▶ C. E. Rasmussen, "Advances in Gaussian Processes," NIPS-2006 Tutorial.