

Gaussian Processes for Learning Nonlinear Functions

Piyush Rai

Topics in Probabilistic Modeling and Inference (CS698X)

Jan 28, 2019



Announcements

- Discussion session on project topics/ideas: Tomorrow 7pm-8pm (KD-101)
- Project proposals due on Feb 1
- HW1 out now. Due on Feb 8, 11:59pm. Please start early.
- Quiz 1 on Jan 31, 7pm-8pm (RM-101)



Recap: Bayesian Generative Classification

- Recall generative classification $p(y = k|x) = \frac{p(y=k)p(x|y=k)}{\sum_{k=1}^K p(y=k)p(x|y=k)}$. Prediction rule for a test input \mathbf{x}_*

$$\begin{aligned} p(y_* = k|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \frac{p(y_* = k|\mathbf{X}, \mathbf{y})p(\mathbf{x}_*|y_* = k, \mathbf{X}, \mathbf{y})}{\sum_{k=1}^K p(y_* = k|\mathbf{X}, \mathbf{y})p(\mathbf{x}_*|y_* = k, \mathbf{X}, \mathbf{y})} \\ &= \frac{p(y_* = k|\mathbf{y})p(\mathbf{x}_*|\mathbf{X}^{(k)})}{\sum_{k=1}^K p(y_* = k|\mathbf{y})p(\mathbf{x}_*|\mathbf{X}^{(k)})} \end{aligned} \quad (1)$$

- Note: $\mathbf{X}^{(k)}$ denotes training inputs with $y = k$
- Here $p(y_* = k|\mathbf{y}) = \int p(y_*|\pi)p(\pi|\mathbf{y})d\pi$ (we did this; recall dice roll example)
- Here $p(\mathbf{x}_*|\mathbf{X}^{(k)}) = \int p(\mathbf{x}_*|\theta_k)p(\theta_k|\mathbf{X}^{(k)})d\theta_k$ (post. predictive dist. of input \mathbf{x}_* under class k)
- Eq (1) is the posterior predictive distribution of test output y_* given input \mathbf{x}_*
 - Note that we have done posterior averaging for all the parameters
- In contrast, for gen. class with MLE/MAP, $p(y_* = k|\mathbf{y}) \approx \pi_k$ and $p(\mathbf{x}_*|\mathbf{X}^{(k)}) \approx p(\mathbf{x}_*|\theta_k)$



Gaussian Processes (GP)

(GP = Bayesian Modeling + Kernel Methods)

(Goal: learning **nonlinear** discriminative models $p(y|\mathbf{x})$)



Linear Models

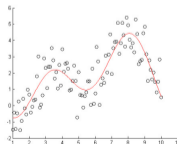
- Consider the problem of learning to map an input $\mathbf{x} \in \mathbb{R}^D$ to an output y
- Linear models use a weighted combination of input features (i.e., $\mathbf{w}^\top \mathbf{x}$) to generate y

$$p(y|\mathbf{w}, \mathbf{x}) = \mathcal{N}(y|\mathbf{w}^\top \mathbf{x}, \beta^{-1}) \quad (\text{Linear Regression})$$

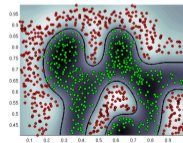
$$p(y|\mathbf{w}, \mathbf{x}) = [\sigma(\mathbf{w}^\top \mathbf{x})]^y [1 - \sigma(\mathbf{w}^\top \mathbf{x})]^{1-y} \quad (\text{Logistic Regression})$$

$$p(y|\mathbf{w}, \mathbf{x}) = \text{ExpFam}(\mathbf{w}^\top \mathbf{x}) \quad (\text{Generalized Linear Model})$$

- The weights \mathbf{w} can be learned using MLE, MAP, or fully Bayesian inference
- However, linear models have limited expressive power. Unable to learn highly nonlinear patterns.



Nonlinear Regression



Nonlinear Classification



Modeling Nonlinear Functions

- Assume the input to output relationship to be modeled by a nonlinear function f

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

$$p(y|f, \mathbf{x}) = [\sigma(f(\mathbf{x}))]^y [1 - \sigma(f(\mathbf{x}))]^{1-y}$$

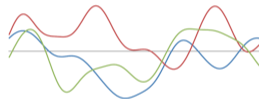
$$p(y|f, \mathbf{x}) = \text{ExpFam}(f(\mathbf{x}))$$

- How can we define such a function nonlinear f ?
- Note: We not only want nonlinearity but also all benefits of probabilistic/Bayesian modeling
 - Must be able to get uncertainty estimates in the function and its predictions
- Usually done in one of the following ways
 - Ad-hoc: Define nonlinear features $\phi(\mathbf{x})$ + train Bayesian linear model ($f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$): HW1)
 - Ad-hoc: Train a neural net to extract features $\phi(\mathbf{x})$ + train Bayesian linear model
 - Bayesian Neural Networks (later this semester)
 - Gaussian Processes (a Bayesian approach to kernel based nonlinear learning; today)



What is Gaussian Process?

- A Gaussian Process, denoted as $\mathcal{GP}(\mu, \kappa)$, defines a **distribution over functions**
 - The GP is defined by **mean function** μ and **covariance/kernel function** κ
- Can use GP as a prior distribution over functions
- Draw from a $\mathcal{GP}(\mu, \kappa)$ will give us a random function f (imagine it as an **infinite dim. vector**)



- Mean function μ models the “average” function f from $\mathcal{GP}(\mu, \kappa)$

$$\mu(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$$

- Cov. function κ models “shape/smoothness” of these functions
 - $\kappa(\cdot, \cdot)$ is a function that computes similarity between two inputs (just like a kernel function)
 - Note: $\kappa(\cdot, \cdot)$ needs to be positive definite (just like kernel functions)
- **Can even learn** μ and especially κ (makes GP very flexible to model, possibly nonlinear, functions)



Gaussian Process

- f is said to be drawn from a $\mathcal{GP}(\mu, \kappa)$ if its finite dim. version is the following joint Gaussian

$$\begin{bmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \\ f(\mathbf{x}_N) \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu(\mathbf{x}_1) \\ \mu(\mathbf{x}_2) \\ \vdots \\ \mu(\mathbf{x}_N) \end{bmatrix}, \begin{bmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_1) \dots \kappa(\mathbf{x}_1, \mathbf{x}_N) \\ \kappa(\mathbf{x}_2, \mathbf{x}_1) \dots \kappa(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots \quad \ddots \quad \vdots \\ \kappa(\mathbf{x}_N, \mathbf{x}_1) \dots \kappa(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix} \right)$$

- The above means that f 's values at any finite set of inputs are jointly Gaussian
- We can also write the above more compactly as $\mathbf{f} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$ where

$$\mathbf{f} = \begin{bmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \\ f(\mathbf{x}_N) \end{bmatrix}, \boldsymbol{\mu} = \begin{bmatrix} \mu(\mathbf{x}_1) \\ \mu(\mathbf{x}_2) \\ \vdots \\ \mu(\mathbf{x}_N) \end{bmatrix}, \mathbf{K} = \begin{bmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_1) \dots \kappa(\mathbf{x}_1, \mathbf{x}_N) \\ \kappa(\mathbf{x}_2, \mathbf{x}_1) \dots \kappa(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots \quad \ddots \quad \vdots \\ \kappa(\mathbf{x}_N, \mathbf{x}_1) \dots \kappa(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

- Note that $p(\mathbf{f}) = \mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$ can be seen as the finite-dimensional version of the GP prior over f
- If mean function is zero, we will have $p(\mathbf{f}) = \mathcal{N}(\mathbf{0}, \mathbf{K})$. **Important:** $p(\mathbf{f}_i | \mathbf{f}_{-i})$ is also Gaussian (where i denotes any subset of inputs and $-i$ denotes rest of the inputs) due to Gaussian properties



A Perspective from Bayesian Linear Regression

- Let's first consider the (probabilistic) linear regression model

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \quad (\text{Prior})$$

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N) \quad (\text{Likelihood w.r.t. } N \text{ obs.})$$

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w} = \mathcal{N}(\mathbf{X}\boldsymbol{\mu}_0, \beta^{-1}\mathbf{I}_N + \mathbf{X}\boldsymbol{\Sigma}_0\mathbf{X}^\top) \quad (\text{Marginal likelihood})$$

$$p(\mathbf{y}|\mathbf{X}) = \mathcal{N}(\mathbf{0}, \beta^{-1}\mathbf{I}_N + \mathbf{X}\mathbf{X}^\top) \quad (\text{if } \boldsymbol{\mu}_0 = \mathbf{0} \text{ and } \boldsymbol{\Sigma}_0 = \mathbf{I})$$

$$p(\mathbf{y}|\mathbf{X}) = \mathcal{N}(\mathbf{0}, \mathbf{X}\mathbf{X}^\top) \quad (\text{if } \beta^{-1} = \infty, \text{ i.e., zero noise})$$

- Thus the **joint marginal distr.** of \mathbf{y} conditioned on \mathbf{X} is the following **multivariate Gaussian**

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{x}_1^\top \mathbf{x}_1 & \dots & \mathbf{x}_1^\top \mathbf{x}_N \\ \mathbf{x}_2^\top \mathbf{x}_1 & \dots & \mathbf{x}_2^\top \mathbf{x}_N \\ \vdots & \ddots & \vdots \\ \mathbf{x}_N^\top \mathbf{x}_1 & \dots & \mathbf{x}_N^\top \mathbf{x}_N \end{bmatrix} \right)$$

- A “function space” view of linear regression as opposed to “weight space” view (both equivalent)



GP for Regression and Classification

(Note that GP only defines the score $f(\mathbf{x})$ but $y = f(\mathbf{x}) + \text{“noise”}$)

(“noise” may be Gaussian, sigmoid-Bernoulli, or something else)



GP Regression

- Training data: $\{\mathbf{x}_n, y_n\}_{n=1}^N$. $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \mathbb{R}$

- Assume the responses to be a noisy function of the inputs

$$y_n = f(\mathbf{x}_n) + \epsilon_n = f_n + \epsilon_n$$

- Assume a zero-mean Gaussian noise: $\epsilon_n \sim \mathcal{N}(\epsilon_n|0, \sigma^2)$

- This implies the following likelihood model: $p(y_n|f_n) = \mathcal{N}(y_n|f_n, \sigma^2)$

- Denote $\mathbf{f} = [f_1, \dots, f_N]$ and $\mathbf{y} = [y_1, \dots, y_N]$. For i.i.d. responses, the joint **likelihood** will be

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

- We now need a **prior** on the function f that enables us to model a nonlinear f

- Let's choose zero mean Gaussian Process prior $\mathcal{GP}(0, \kappa)$ on f , which is equivalent to

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

where $K_{nm} = \kappa(\mathbf{x}_n, \mathbf{x}_m)$. For now, assume κ is a known function with fixed hyperparameters.



GP Regression

- The likelihood model: $p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2\mathbf{I}_N)$. The prior distribution: $p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$
- The posterior $p(\mathbf{f}|\mathbf{y}) \propto p(\mathbf{f})p(\mathbf{y}|\mathbf{f})$, which will be another Gaussian (Exercise: Find its expression)
- What's the posterior predictive $p(y_*|\mathbf{x}_*, \mathbf{y}, \mathbf{X})$ or $p(y_*|\mathbf{y})$ (skipping \mathbf{X}, \mathbf{x}_* from the notation)?

$$p(y_*|\mathbf{y}) = \int p(y_*|f_*)p(f_*|\mathbf{y})df_*$$

where $p(f_*|\mathbf{y}) = \int p(f_*|\mathbf{f})p(\mathbf{f}|\mathbf{y})d\mathbf{f}$ and note that $p(f_*|\mathbf{f})$ must be Gaussian for GP

- For this case (GP regression), we actually don't need to compute $p(y_*|\mathbf{y})$ using the above method
- Reason: The **marginal distribution** of the training data responses \mathbf{y}

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K} + \sigma^2\mathbf{I}_N) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C}_N)$$

- Using the same result, the marginal distribution $p(y_*) = \mathcal{N}(y_*|0, \kappa(\mathbf{x}_*, \mathbf{x}_*) + \sigma^2)$



GP Regression: Making Predictions

- Let's consider the joint distr. of N training responses \mathbf{y} and test response y_*

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \mathbf{C}_{N+1}\right)$$

where the $(N+1) \times (N+1)$ matrix \mathbf{C}_{N+1} is given by

$$\mathbf{C}_{N+1} = \begin{bmatrix} \mathbf{C}_N & \mathbf{k}_* \\ \mathbf{k}_*^\top & c \end{bmatrix}$$

and $\mathbf{k}_* = [\kappa(\mathbf{x}_*, \mathbf{x}_1), \dots, \kappa(\mathbf{x}_*, \mathbf{x}_N)]^\top$, $c = \kappa(\mathbf{x}_*, \mathbf{x}_*) + \sigma^2$

- The desired **predictive posterior** will be (using conditional from joint property of Gaussian)

$$\begin{aligned} p(y_* | \mathbf{y}) &= \mathcal{N}(y_* | \mu_*, \sigma_*^2) \\ \mu_* &= \mathbf{k}_*^\top \mathbf{C}_N^{-1} \mathbf{y} \\ \sigma_*^2 &= \kappa(\mathbf{x}_*, \mathbf{x}_*) + \sigma^2 - \mathbf{k}_*^\top \mathbf{C}_N^{-1} \mathbf{k}_* \end{aligned}$$



GP Regression: Interpreting GP Predictions

- Let's look at the predictions made by GP regression

$$\begin{aligned}p(y_*|\mathbf{y}) &= \mathcal{N}(y_*|\mu_*, \sigma_*^2) \\ \mu_* &= \mathbf{k}_*^\top \mathbf{C}_N^{-1} \mathbf{y} \\ \sigma_*^2 &= k(\mathbf{x}_*, \mathbf{x}_*) + \sigma^2 - \mathbf{k}_*^\top \mathbf{C}_N^{-1} \mathbf{k}_*\end{aligned}$$

- Two interpretations for the mean prediction μ_*

- A kernel SVM like interpretation

$$\mu_* = \mathbf{k}_*^\top \mathbf{C}_N^{-1} \mathbf{y} = \mathbf{k}_*^\top \boldsymbol{\alpha} = \sum_{n=1}^N k(\mathbf{x}_*, \mathbf{x}_n) \alpha_n$$

where $\boldsymbol{\alpha}$ is akin to the weights of support vectors

- A nearest neighbors interpretation

$$\mu_* = \mathbf{k}_*^\top \mathbf{C}_N^{-1} \mathbf{y} = \mathbf{w}^\top \mathbf{y} = \sum_{n=1}^N w_n y_n$$

where \mathbf{w} is akin to the weights of the neighbors



Next Class

- Properties of GP based models, choice of kernels, etc
- Learning hyperparameters in GP based models
- GP for classification and GLMs
- Making GP models scalable
- Some recent advances in GP

