

# Gaussian Processes (Contd.), Intro to Latent Variable Models

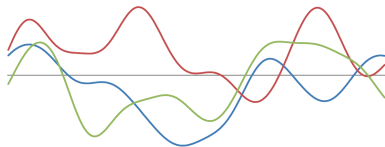
Piyush Rai

Topics in Probabilistic Modeling and Inference (CS698X)

Feb 6, 2018

# Gaussian Process

- Denoted as  $\mathcal{GP}(\mu, \kappa)$ . It is a **distribution over functions**
  - Defined by **mean function**  $\mu$  and **covariance function**  $\kappa$



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

$$\mu(\mathbf{x}) = \mathbb{E}_{f \sim \mathcal{GP}(\mu, \kappa)}[f(\mathbf{x})]$$

- Cov. function  $\kappa$  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)
- An appealing aspect: Can learn “hyperparameters”  $\mu$  and  $\kappa$  (i.e., can learn the kernel)

# Gaussian Process

- A random function  $f$  drawn from  $\mathcal{GP}(\mu, \kappa)$ , has the following “finite dimensional marginal”

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

- Note that, unlike traditional parametric definitions  $y_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n$ , here the “function”  $f$  is defined in terms of the **joint marginal distribution** of a set of outputs given the corresponding inputs
  - **Weight-space view** vs **function-space view** of input to output mapping
- We can also write the above more compactly as  $\mathbf{f} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$  where

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

Note:  $\mathbf{K}$  is also called the **kernel matrix**.  $K_{nm} = \kappa(\mathbf{x}_n, \mathbf{x}_m)$

- Therefore the finite-dim version of the GP  $p(f) = \mathcal{GP}(\mu, \kappa)$  is  $p(\mathbf{f}) = \mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$

# Using GP as a Prior Distribution in Regression/Classification

- Consider a problem with data of the form  $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$  using a discriminative model
- Typically the probability distribution of  $y_n$  given  $\mathbf{x}_n$  depends on a “score”  $f_n = f(\mathbf{x}_n)$ , e.g.,

$$\begin{array}{ll} \text{Regression} & p(y_n|\mathbf{x}_n) = p(y_n|f_n) = \mathcal{N}(y_n|f_n, \sigma^2) \\ \text{Classification} & p(y_n|\mathbf{x}_n) = p(y_n|f_n) = \text{Bernoulli}(y_n|\sigma(f_n)) \end{array}$$

- Suppose the vector of scores is defined as  $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]$
- Can use GP as a prior on the function  $f$ , or equivalently on the score vector  $\mathbf{f}$

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{0}, \mathbf{K}) \quad (\text{assuming zero mean function})$$

- Can now combine this prior with the likelihood  $p(\mathbf{y}|\mathbf{f})$  to compute
  - Posterior over the function  $f$ , which is given in the form  $p(\mathbf{f}|\mathbf{y}) \propto p(\mathbf{f})p(\mathbf{y}|\mathbf{f})$
  - Posterior predictive  $p(y_*|\mathbf{y}) = \int p(y_*|f_*)p(f_*|\mathbf{y})df_*$  where  $p(f_*|\mathbf{y}) = \int \underbrace{p(f_*|\mathbf{f})}_{\text{Always Gaussian}} p(\mathbf{f}|\mathbf{y})d\mathbf{f}$

Always  
Gaussian

# Scalability Aspects of GP

- Computational costs in some steps of GP based models scale in the size of training data
  - E.g., test time prediction in GP regression takes  $O(N)$  time

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

- GP models often require matrix inversions - takes  $O(N^3)$  time. Storage also requires  $O(N^2)$  space
- A lot of work on speeding up GPs. For some recent advances, you may check out
  - “Thoughts on Massively Scalable Gaussian Processes”
  - “Kernel Interpolation for Scalable Structured Gaussian Processes (KISS-GP)”
  - .. and references therein
- Note that nearest neighbor methods and kernel methods also face similar issues w.r.t. scalability
  - Many tricks to speed up kernel methods can be used for speeding up GPs too

# GP: A few comments

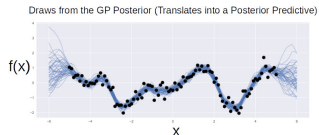
- GP is a **nonparametric model**. Why called “nonparametric”?
  - Complexity (representation size) of the function  $f$  grows in the size of training data
  - To see this, note the form of the GP predictions, e.g., predictive mean in GP regression

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

- It implies that  $f(\cdot) = \sum_{n=1}^N \alpha_n k(\cdot, \mathbf{x}_n)$ , which means  $f$  is written in terms of all training examples
  - Thus the representation size of  $f$  depends on the number of training examples
- In contrast, a parametric model has a size that doesn't grow with training data
  - E.g., a linear model learns a fixed-sized weight vector  $\mathbf{w} \in \mathbb{R}^D$  ( $D$  parameters, size independent of  $N$ )
- Nonparametric models therefore are more flexible since their complexity is not limited beforehand
  - Note: Methods such as nearest neighbors and kernel SVMs are also nonparametric (but not Bayesian)
- GPs equivalent to **infinitely-wide single hidden-layer neural net** (under some technical conditions)

# GP: A few other comments

- Can be thought of as Bayesian analogues of kernel methods
  - Can get estimate in the uncertainty in the function and its predictions



- Can learn the kernel (by learning the hyperparameters of the kernels)
- Not limited to supervised learning problems
  - The function  $f$  could even be a mapping of an unknown quantity to an observed quantity

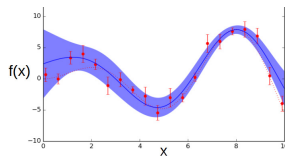
$$x_n = f(z_n) + \text{"noise"}$$

where  $z_n$  is a latent representation of  $x_n$  (**"GP latent variable models"** for nonlin. dim. red.)

- Many mature implementations of GP exist. You may check out
  - GPML (MATLAB), GPsuff (MATLAB/Octave), GPy (Python), GPyTorch (PyTorch)

# GPs are very versatile!

- GPs enable us to learn nonlinear functions while also capturing the uncertainty



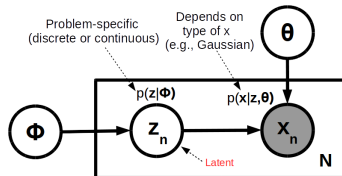
- Uncertainty can tell us where to acquire more training data to improve the function's estimate
  - Especially useful if we can't get too many training examples (e.g., expensive inputs and/or labels)
- This is very useful in a wide range of applications involving sequential decision-making
  - **Active Learning**: Learning a function by gathering the most informative training examples
  - **Bayesian Optimization**: Optimizing an expensive to evaluate functions (and maybe we don't even know its form) – boils down to simultaneous function learning and optimization
- We will look at some of these later during the semester



# Latent Variable Models

# Modeling Data via Latent Variables

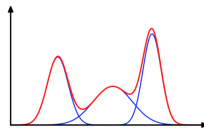
- Assume each observation  $x_n$  to be associated with a latent variable  $z_n$



- $z_n$  is a latent representation or “encoding” of  $x_n$  (often reveals the **latent structure** of data)
  - $z_n \in \{1, \dots, K\}$  denotes the cluster  $x_n$  belongs to
  - $z_n \in \mathbb{R}^K$  denotes a low-dimensional latent representation or latent “code” for  $x_n$
- Sometimes  $\mathbf{Z} = \{z_1, \dots, z_N\}$  called **“local” variables**;  $\Theta = (\theta, \phi)$  called **“global” variables**
  - The goal is to learn these unknowns given the observed data  $\mathbf{X} = \{x_1, \dots, x_N\}$ , i.e.,  $p(\mathbf{Z}, \Theta | \mathbf{X})$

# Motivating Example 1: Mixture Model

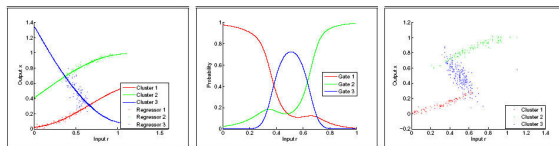
- Assume data  $\{\mathbf{x}_n\}_{n=1}^N$  was generated from a mixture of  $K$  distributions  $p(\mathbf{x}|\theta_1), \dots, p(\mathbf{x}|\theta_K)$



- Don't know which component generated each  $\mathbf{x}_n$  (o/w it is simply **generative classification**)
- Assume a latent random variable  $\mathbf{z}_n \in \{1, \dots, K\}$  denotes which component generated  $\mathbf{x}_n$
- Here is a simple **generative story** for each  $\mathbf{x}_n$ ,  $n = 1, 2, \dots, N$ 
  - First choose a mixture component  $\mathbf{z}_n \in \{1, 2, \dots, K\}$  as  $\mathbf{z}_n \sim \text{multinoulli}(\mathbf{z}|\pi)$
  - Now generate  $\mathbf{x}_n$  from that mixture component as  $\mathbf{x}_n \sim p(\mathbf{x}|\theta_{\mathbf{z}_n})$
- Goal: Given data  $\{\mathbf{x}_n\}_{n=1}^N$ , learn the  $K$  distributions  $(\theta_1, \dots, \theta_K)$  and latent variables  $\mathbf{z}_1, \dots, \mathbf{z}_N$
- If each  $p(\mathbf{x}|\theta_k)$  is a Gaussian  $\Rightarrow$  **Gaussian Mixture Model** (used for probabilistic or “soft” clustering)

# Motivating Example 2: Mixture of Experts

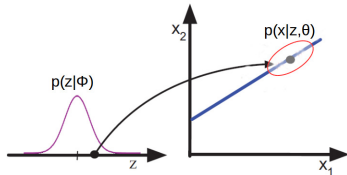
- Assume  $N$  inputs  $\mathbf{x}_1, \dots, \mathbf{x}_N$  whose responses are  $y_1, \dots, y_N$
- Assume responses  $y_1, \dots, y_N$  are generated by a mixture of  $K$  linear regression models (experts)



- Each linear regression model handles one region of the input space (regions may have overlap)
- The resulting  $\mathbf{x}$  to  $y$  mapping learned by such a mixture model is effectively a nonlinear model
- The **generative story** for each  $y_n$  conditioned on  $\mathbf{x}_n$ 
  - First choose an expert  $z_n \in \{1, 2, \dots, K\}$  as  $z_n \sim \text{multinoulli}(z | \pi(\mathbf{x}_n))$  (note:  $\pi$  depends on  $\mathbf{x}_n$  too)
  - Now generate  $y_n$  using the regression model of component  $K$  as  $y_n | \mathbf{x}_n \sim p(y_n | \mathbf{x}_n, \mathbf{w}_{z_n})$

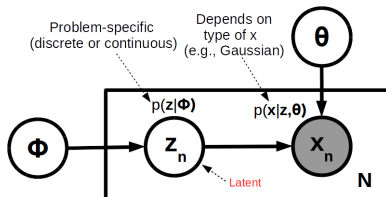
# Motivating Example 3: Latent Factor Model

- Assume data  $\mathbf{x}_n \in \mathbb{R}^D$  generated from a low-dimensional latent factor  $\mathbf{z}_n \in \mathbb{R}^K$



- The  $\mathbf{z}$  to  $\mathbf{x}$  map can be a linear/nonlinear transformation
- Consider the following **generative story** for each  $\mathbf{x}_n$ ,  $n = 1, 2, \dots, N$ 
  - First generate  $\mathbf{z}_n$  from a  $K$ -dim distr. as  $\mathbf{z}_n \sim p(\mathbf{z}|\phi)$
  - Now generate  $\mathbf{x}_n$  from a  $D$ -dim distr. as  $\mathbf{x}_n \sim p(\mathbf{x}|\mathbf{z}_n, \theta)$
- If  $p(\mathbf{z}|\phi)$  and  $p(\mathbf{x}|\mathbf{z}, \theta)$  are Gaussians and  $\mathbf{z}$  to  $\mathbf{x}$  map linear  $\Rightarrow$  **factor analysis** or **probabilistic PCA**

# Learning Parameters of Latent Variable Models



- Suppose we want to learn the “global” parameters  $\Theta = (\theta, \phi)$  of this LVM
- While full posterior inference can be tried for  $\Theta$ , let’s keep it simple for now and do MLE/MAP
- The MLE would be for the above LVM would be

$$\hat{\Theta} = \arg \max_{\Theta} \sum_{n=1}^N \log p(\mathbf{x}_n | \Theta)$$

- It turns out that the above problem is a hard problem in general!

# Why is MLE/MAP Hard for LVMs?

- First of all, note that when working with exp. family distributions, MLE turns out to be very easy
  - Reason: Usually **log of an exp-fam distribution** has simple algebraic form, easy to do MLE
- In the MLE/MAP for LVM, the goal is to solve  $\hat{\Theta} = \arg \max_{\Theta} \sum_{n=1}^N \log p(\mathbf{x}_n | \Theta)$
- To do so, we first need  $p(\mathbf{x}_n | \Theta)$  which is a marginal

$$\text{Discrete } \mathbf{z}_n: \quad p(\mathbf{x}_n | \Theta) = \sum p(\mathbf{x}_n, \mathbf{z}_n | \Theta) = \sum_{k=1}^K p(\mathbf{x}_n | \mathbf{z}_n = k, \Theta) p(\mathbf{z}_n = k | \Theta)$$

$$\text{Continuous } \mathbf{z}_n: \quad p(\mathbf{x}_n | \Theta) = \int p(\mathbf{x}_n, \mathbf{z}_n | \Theta) d\mathbf{z}_n = \int p(\mathbf{x}_n | \mathbf{z}_n, \Theta) p(\mathbf{z}_n | \Theta) d\mathbf{z}_n$$

- **Causes problem!** Even if  $p(\mathbf{x}_n | \mathbf{z}_n, \Theta)$  and  $p(\mathbf{z}_n | \Theta)$  (and thus the joint  $p(\mathbf{x}_n, \mathbf{z}_n | \Theta)$ ) are exp-fam distributions, marginal  $p(\mathbf{x}_n | \Theta)$  isn't in general an exp-fam distribution (Bar-Lev et al, 1994)
  - Consequently,  $\log p(\mathbf{x}_n | \Theta)$  doesn't usually have a simple algebraic form amenable for MLE
- On the other hand, if someone gave us the “good guess”  $\hat{\mathbf{z}}_n$  of  $\mathbf{z}_n$ , MLE becomes much simpler (that's the underlying idea behind the Expectation Maximization algorithm)