

Learning Nonlinear Functions via Gaussian Processes (1)

Piyush Rai

Probabilistic Machine Learning (CS772A)

October 17, 2017

Linear Models

- Consider the problem of learning to map an input $\mathbf{x} \in \mathbb{R}^D$ to an output y

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

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

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

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

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

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

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

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

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

- This can also be extended to [generalized linear models](#) (to model other types of y)

$$p(y|\mathbf{w}, \mathbf{x}) = \text{GLM}(y|\eta) \quad \text{where } \eta = \mathbf{w}^\top \mathbf{x}$$

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

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

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

- This can also be extended to [generalized linear models](#) (to model other types of y)

$$p(y|\mathbf{w}, \mathbf{x}) = \text{GLM}(y|\eta) \quad \text{where } \eta = \mathbf{w}^\top \mathbf{x}$$

- The weights \mathbf{w} can be learned using MLE, MAP, or fully Bayesian inference

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

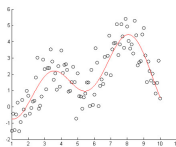
$$p(y|\mathbf{w}, \mathbf{x}) = \mathcal{N}(y|\mathbf{w}^\top \mathbf{x}, \beta^{-1}) \quad (\text{Linear Regression})$$

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

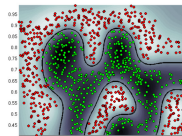
- This can also be extended to **generalized linear models** (to model other types of y)

$$p(y|\mathbf{w}, \mathbf{x}) = \text{GLM}(y|\eta) \quad \text{where } \eta = \mathbf{w}^\top \mathbf{x}$$

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



Nonlinear Regression



Nonlinear Classification

Linear Regression

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

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

$$p(y_n | \mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) \quad (\text{Likelihood w.r.t. one obs.})$$

Linear Regression

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

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

$$p(y_n | \mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) \quad (\text{Likelihood w.r.t. one obs.})$$

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

Linear Regression

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

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

$$p(y_n|\mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) \quad (\text{Likelihood w.r.t. one obs.})$$

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

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

Linear Regression

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

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

$$p(y_n|\mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) \quad (\text{Likelihood w.r.t. one obs.})$$

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

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

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

Linear Regression

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

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

$$p(y_n|\mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) \quad (\text{Likelihood w.r.t. one obs.})$$

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

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

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

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

Linear Regression

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

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

$$p(y_n|\mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) \quad (\text{Likelihood w.r.t. one obs.})$$

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

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

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

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

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

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

Linear Regression

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

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

$$p(y_n|\mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) \quad (\text{Likelihood w.r.t. one obs.})$$

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

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

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

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

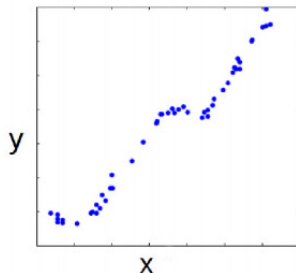
- Thus the **joint marginal distr.** of \mathbf{y} conditioned on \mathbf{X} is the following **multivariate normal**

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

- Note that the cov. matrix above tells us how the different outputs “co-vary” w.r.t. each other

Nonlinear Regression

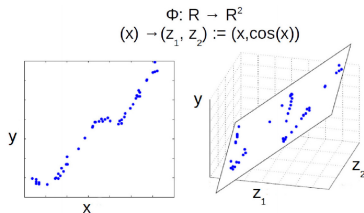
- Consider the following dataset for regression



- Input-output relationship is NOT linear
- A linear model would do badly on such data!

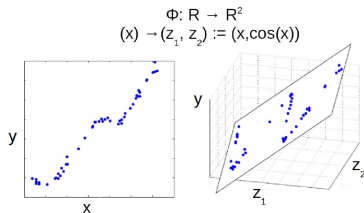
Nonlinear Regression via Mapping Inputs to Higher Dimensions

- Suppose we replace each input \mathbf{x}_n by a higher dimensional mapping $\phi(\mathbf{x}_n)$



Nonlinear Regression via Mapping Inputs to Higher Dimensions

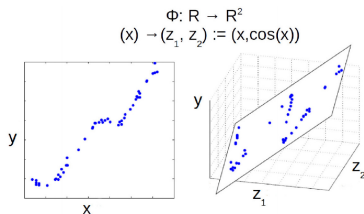
- Suppose we replace each input \mathbf{x}_n by a higher dimensional mapping $\phi(\mathbf{x}_n)$



- Now the nonlinear regression problem becomes a linear reg. problem in a higher dim. space

Nonlinear Regression via Mapping Inputs to Higher Dimensions

- Suppose we replace each input \mathbf{x}_n by a higher dimensional mapping $\phi(\mathbf{x}_n)$

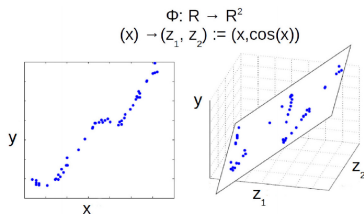


- Now the nonlinear regression problem becomes a linear reg. problem in a higher dim. space
- The joint distribution \mathbf{y} given \mathbf{X} will be the same as the linear case, except that each entry $\mathbf{x}_n^\top \mathbf{x}_m$ in the covariance matrix is replaced by $\phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$

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

Nonlinear Regression via Mapping Inputs to Higher Dimensions

- Suppose we replace each input \mathbf{x}_n by a higher dimensional mapping $\phi(\mathbf{x}_n)$



- Now the nonlinear regression problem becomes a linear reg. problem in a higher dim. space
- The joint distribution \mathbf{y} given \mathbf{X} will be the same as the linear case, except that each entry $\mathbf{x}_n^\top \mathbf{x}_m$ in the covariance matrix is replaced by $\phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$

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

- Note that $\phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$ is the Euclidean similarity between \mathbf{x}_n and \mathbf{x}_m in the high-dim space

A Notion of Similarity: Kernel Functions

- Consider a covariance/kernel function κ that measures similarity between two inputs

A Notion of Similarity: Kernel Functions

- Consider a covariance/kernel function κ that measures similarity between two inputs
- Suppose κ corresponds to **implicitly mapping data** to a higher dimensional space via a feature mapping ϕ ($\mathbf{x} \rightarrow \phi(\mathbf{x})$) and computing the dot product (Euclidean) similarity in that space

$$\kappa(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$$

A Notion of Similarity: Kernel Functions

- Consider a covariance/kernel function κ that measures similarity between two inputs
- Suppose κ corresponds to **implicitly mapping data** to a higher dimensional space via a feature mapping ϕ ($\mathbf{x} \rightarrow \phi(\mathbf{x})$) and computing the dot product (Euclidean) similarity in that space

$$\kappa(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$$

- An example: Radial basis function (RBF) kernel: $\kappa(\mathbf{x}_n, \mathbf{x}_m) = \exp\left(-\frac{\|\mathbf{x}_n - \mathbf{x}_m\|^2}{\gamma}\right)$ where ϕ_n corresponds to an infinite dimensional mapping. Another example: Polynomial kernel

A Notion of Similarity: Kernel Functions

- Consider a covariance/kernel function κ that measures similarity between two inputs
- Suppose κ corresponds to **implicitly mapping data** to a higher dimensional space via a feature mapping ϕ ($\mathbf{x} \rightarrow \phi(\mathbf{x})$) and computing the dot product (Euclidean) similarity in that space

$$\kappa(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$$

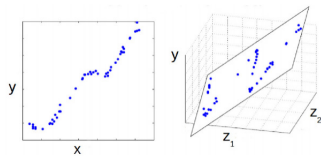
- An example: Radial basis function (RBF) kernel: $\kappa(\mathbf{x}_n, \mathbf{x}_m) = \exp\left(-\frac{\|\mathbf{x}_n - \mathbf{x}_m\|^2}{\gamma}\right)$ where ϕ_n corresponds to an infinite dimensional mapping. Another example: Polynomial kernel
- Allows extending linear models to nonlinear models without explicitly computing $\phi(\mathbf{x}_n)$

A Notion of Similarity: Kernel Functions

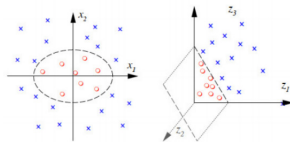
- Consider a covariance/kernel function κ that measures similarity between two inputs
- Suppose κ corresponds to **implicitly mapping data** to a higher dimensional space via a feature mapping ϕ ($\mathbf{x} \rightarrow \phi(\mathbf{x})$) and computing the dot product (Euclidean) similarity in that space

$$\kappa(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$$

- An example: Radial basis function (RBF) kernel: $\kappa(\mathbf{x}_n, \mathbf{x}_m) = \exp\left(-\frac{\|\mathbf{x}_n - \mathbf{x}_m\|^2}{\gamma}\right)$ where ϕ_n corresponds to an infinite dimensional mapping. Another example: Polynomial kernel
- Allows extending linear models to nonlinear models without explicitly computing $\phi(\mathbf{x}_n)$



Nonlinear Regression



Nonlinear Classification

Nonlinear Regression

- We saw that, for a nonlin. reg. model with each \mathbf{x}_n replaced by a high-dim map $\phi(\mathbf{x}_n)$

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

Nonlinear Regression

- We saw that, for a nonlin. reg. model with each \mathbf{x}_n replaced by a high-dim map $\phi(\mathbf{x}_n)$

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

- Writing $\kappa(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$, we have

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

Nonlinear Regression

- We saw that, for a nonlin. reg. model with each \mathbf{x}_n replaced by a high-dim map $\phi(\mathbf{x}_n)$

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

- Writing $\kappa(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$, we have

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

- More compactly, we can write $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$ where $C_{nm} = \kappa(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$

Nonlinear Regression

- We saw that, for a nonlin. reg. model with each \mathbf{x}_n replaced by a high-dim map $\phi(\mathbf{x}_n)$

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

- Writing $\kappa(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$, we have

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

- More compactly, we can write $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$ where $C_{nm} = \kappa(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m)$
- Note: In the noisy case, we only need to add β^{-1} ($=\sigma^2$) to diagonals of \mathbf{C}

Nonlinear Regression

- Assume noisy case. Consider joint distr. of N training outputs \mathbf{y} and test output y_* (input \mathbf{x}_*)

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

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

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

Nonlinear Regression

- Assume noisy case. Consider joint distr. of N training outputs \mathbf{y} and test output y_* (input \mathbf{x}_*)

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

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

Nonlinear Regression

- Assume noisy case. Consider joint distr. of N training outputs \mathbf{y} and test output y_* (input \mathbf{x}_*)

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

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

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

Nonlinear Regression

- Given the jointly Gaussian distribution

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{C}_N & \mathbf{k}_* \\ \mathbf{k}_*^\top & c \end{bmatrix}\right)$$

Nonlinear Regression

- Given the jointly Gaussian distribution

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{C}_N & \mathbf{k}_* \\ \mathbf{k}_*^\top & c \end{bmatrix}\right)$$

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

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

Nonlinear Regression

- Given the jointly Gaussian distribution

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{C}_N & \mathbf{k}_* \\ \mathbf{k}_*^\top & c \end{bmatrix}\right)$$

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

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

- How does this compare with linear model's posterior predictive?

Nonlinear Regression

- Given the jointly Gaussian distribution

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{C}_N & \mathbf{k}_* \\ \mathbf{k}_*^\top & c \end{bmatrix}\right)$$

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

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

- How does this compare with linear model's posterior predictive?
 - One key difference is that the model's prediction now **depends explicitly on all the training data**

Nonlinear Regression

- Given the jointly Gaussian distribution

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{C}_N & \mathbf{k}_* \\ \mathbf{k}_*^\top & c \end{bmatrix}\right)$$

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

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

- How does this compare with linear model's posterior predictive?
 - One key difference is that the model's prediction now **depends explicitly on all the training data**
 - Such models are called **"nonparametric" models** (model's complexity depends on the data)

Nonlinear Regression

- Given the jointly Gaussian distribution

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{C}_N & \mathbf{k}_* \\ \mathbf{k}_*^\top & c \end{bmatrix}\right)$$

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

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

- How does this compare with linear model's posterior predictive?
 - One key difference is that the model's prediction now **depends explicitly on all the training data**
 - Such models are called **"nonparametric" models** (model's complexity depends on the data)
 - A downside: Making predictions gets slower

Learning Nonlinear Functions via Gaussian Processes (what we just saw, but in a more formal way)

Gaussian Process

- Recall that our kernel version of probabilistic linear regression (noiseless case) had

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

Gaussian Process

- Recall that our kernel version of probabilistic linear regression (noiseless case) had

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

- Note: A slightly general version would be with nonzero mean

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

where μ is a “mean function”

Gaussian Process

- Recall that our kernel version of probabilistic linear regression (noiseless case) had

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

- Note: A slightly general version would be with nonzero mean

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

where μ is a “mean function”

- Suppose we define $y_n = f(\mathbf{x}_n)$ where f is the (possibly nonlinear) function that maps \mathbf{x} to y

Gaussian Process

- Recall that our kernel version of probabilistic linear regression (noiseless case) had

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

- Note: A slightly general version would be with nonzero mean

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

where μ is a “mean function”

- Suppose we define $y_n = f(\mathbf{x}_n)$ where f is the (possibly nonlinear) function that maps \mathbf{x} to y
- Gaussian Process (GP) formally defines a **distribution over nonlinear functions** and help us perform Bayesian inference on such functions

Gaussian Process

- Recall that our kernel version of probabilistic linear regression (noiseless case) had

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

- Note: A slightly general version would be with nonzero mean

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

where μ is a “mean function”

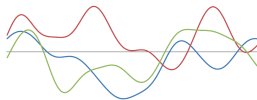
- Suppose we define $y_n = f(\mathbf{x}_n)$ where f is the (possibly nonlinear) function that maps \mathbf{x} to y
- Gaussian Process (GP) formally defines a **distribution over nonlinear functions** and help us perform Bayesian inference on such functions (important: not limited to only regression problems)

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 function** κ

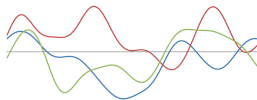
Gaussian Process

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



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 function** κ
- Draw from a $\mathcal{GP}(\mu, \kappa)$ will give us a random function f (imagine it as an **infinite dim. vector**)

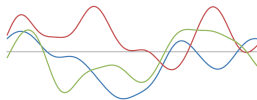


- Mean function μ models the “average” function f from $\mathcal{GP}(\mu, \kappa)$

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

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 function** κ
- Draw from a $\mathcal{GP}(\mu, \kappa)$ will give us a random function f (imagine it as an **infinite dim. vector**)



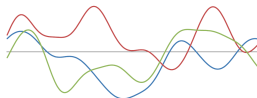
- Mean function μ models the “average” function f from $\mathcal{GP}(\mu, \kappa)$

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

- Cov. function κ models “shape/smoothness” of these functions

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 function** κ
- Draw from a $\mathcal{GP}(\mu, \kappa)$ will give us a random function f (imagine it as an **infinite dim. vector**)



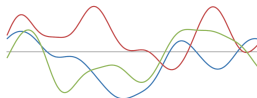
- Mean function μ models the “average” function f from $\mathcal{GP}(\mu, \kappa)$

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

- Cov. function κ models “shape/smoothness” of these functions
 - $\kappa(.,.)$ is a function that computes similarity between two inputs (just like a kernel function)

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 function** κ
- Draw from a $\mathcal{GP}(\mu, \kappa)$ will give us a random function f (imagine it as an **infinite dim. vector**)



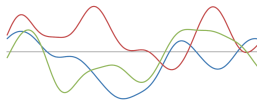
- Mean function μ models the “average” function f from $\mathcal{GP}(\mu, \kappa)$

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

- Cov. function κ models “shape/smoothness” of these functions
 - $\kappa(.,.)$ is a function that computes similarity between two inputs (just like a kernel function)
 - Note: $\kappa(.,.)$ needs to be positive definite (just like kernel functions)

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 function** κ
- Draw from a $\mathcal{GP}(\mu, \kappa)$ will give us a random function f (imagine it as an **infinite dim. vector**)



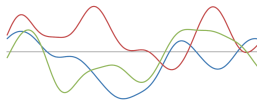
- Mean function μ models the “average” function f from $\mathcal{GP}(\mu, \kappa)$

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

- Cov. function κ models “shape/smoothness” of these functions
 - $\kappa(.,.)$ is a function that computes similarity between two inputs (just like a kernel function)
 - Note: $\kappa(.,.)$ 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

- A Gaussian Process, denoted as $\mathcal{GP}(\mu, \kappa)$, defines a **distribution over functions**
 - The GP is defined by **mean function** μ and **covariance function** κ
- 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
 - $\kappa(.,.)$ is a function that computes similarity between two inputs (just like a kernel function)
 - Note: $\kappa(.,.)$ 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)
- GP can therefore be used as a flexible **prior distribution** over functions

Gaussian Process Prior

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

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

Gaussian Process Prior

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

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

- The above means that f 's output at any finite set of inputs is jointly Gaussian

Gaussian Process Prior

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

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

- The above means that f 's output at any finite set of inputs is jointly Gaussian
- Also makes intuitive sense: If $k(\mathbf{x}_n, \mathbf{x}_m)$ is large, we would expect $f(\mathbf{x}_n)$ and $f(\mathbf{x}_m)$ be the close

Gaussian Process Prior

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

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

- The above means that f 's output at any finite set of inputs is jointly Gaussian
- Also makes intuitive sense: If $k(\mathbf{x}_n, \mathbf{x}_m)$ is large, we would expect $f(\mathbf{x}_n)$ and $f(\mathbf{x}_m)$ be the close
- We can also write the above more compactly as $\mathbf{f} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$ where

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

Gaussian Process Prior

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

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

- The above means that f 's output at any finite set of inputs is jointly Gaussian
- Also makes intuitive sense: If $k(\mathbf{x}_n, \mathbf{x}_m)$ is large, we would expect $f(\mathbf{x}_n)$ and $f(\mathbf{x}_m)$ be the close
- 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**. $K_{nm} = \kappa(\mathbf{x}_n, \mathbf{x}_m)$

Gaussian Process Prior

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

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

- The above means that f 's output at any finite set of inputs is jointly Gaussian
- Also makes intuitive sense: If $k(\mathbf{x}_n, \mathbf{x}_m)$ is large, we would expect $f(\mathbf{x}_n)$ and $f(\mathbf{x}_m)$ be the close
- 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**. $K_{nm} = \kappa(\mathbf{x}_n, \mathbf{x}_m)$

- 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**

Gaussian Process Prior

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

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

- The above means that f 's output at any finite set of inputs is jointly Gaussian
- Also makes intuitive sense: If $k(\mathbf{x}_n, \mathbf{x}_m)$ is large, we would expect $f(\mathbf{x}_n)$ and $f(\mathbf{x}_m)$ be the close
- 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**. $K_{nm} = \kappa(\mathbf{x}_n, \mathbf{x}_m)$

- 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 the mean function is zero, we will have $p(\mathbf{f}) = \mathcal{N}(\mathbf{0}, \mathbf{K})$

Gaussian Process Regression

GP Regression

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

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

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

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: $\epsilon_n \sim \mathcal{N}(\epsilon_n|0, \sigma^2)$

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: $p(y_n|f_n) = \mathcal{N}(y_n|f_n, \sigma^2)$

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 $\mathbf{y} = [y_1, \dots, y_N]$.

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

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

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

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 $K_{nm} = \kappa(\mathbf{x}_n, \mathbf{x}_m)$.

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

GP Regression

- The likelihood model: $p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I}_N)$
- The prior distribution: $p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$

GP Regression

- The likelihood model: $p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I}_N)$
- The prior distribution: $p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$
- Note: If we only care about making predictions on new test inputs, we don't actually need to compute the posterior $p(\mathbf{f}|\mathbf{y})$ here (although, due to Gaussians, we can do so easily)

GP Regression

- The likelihood model: $p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2\mathbf{I}_N)$
- The prior distribution: $p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$
- Note: If we only care about making predictions on new test inputs, we don't actually need to compute the posterior $p(\mathbf{f}|\mathbf{y})$ here (although, due to Gaussians, we can do so easily)
- The **marginal distribution** of the training data responses \mathbf{y} (after integrating out the “unknown” \mathbf{f})

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

GP Regression

- The likelihood model: $p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2\mathbf{I}_N)$
- The prior distribution: $p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$
- Note: If we only care about making predictions on new test inputs, we don't actually need to compute the posterior $p(\mathbf{f}|\mathbf{y})$ here (although, due to Gaussians, we can do so easily)
- The **marginal distribution** of the training data responses \mathbf{y} (after integrating out the “unknown” \mathbf{f})

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

- What will be the prediction y_* for a new test input \mathbf{x}_* ?

GP Regression

- The likelihood model: $p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2\mathbf{I}_N)$
- The prior distribution: $p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$
- Note: If we only care about making predictions on new test inputs, we don't actually need to compute the posterior $p(\mathbf{f}|\mathbf{y})$ here (although, due to Gaussians, we can do so easily)
- The **marginal distribution** of the training data responses \mathbf{y} (after integrating out the “unknown” \mathbf{f})

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

- What will be the prediction y_* for a new test input \mathbf{x}_* ?
- Well, using the above result, we know that the marginal distribution of y_* too must be

$$p(y_*) = \mathcal{N}(y_*|0, \kappa(\mathbf{x}_*, \mathbf{x}_*) + \sigma^2)$$

GP Regression

- The likelihood model: $p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2\mathbf{I}_N)$
- The prior distribution: $p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$
- Note: If we only care about making predictions on new test inputs, we don't actually need to compute the posterior $p(\mathbf{f}|\mathbf{y})$ here (although, due to Gaussians, we can do so easily)
- The **marginal distribution** of the training data responses \mathbf{y} (after integrating out the “unknown” \mathbf{f})

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

- What will be the prediction y_* for a new test input \mathbf{x}_* ?
- Well, using the above result, we know that the marginal distribution of y_* too must be

$$p(y_*) = \mathcal{N}(y_*|0, \kappa(\mathbf{x}_*, \mathbf{x}_*) + \sigma^2)$$

- But what we actually want is the **predictive posterior** $p(y_*|\mathbf{y})$

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

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

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

GP Regression: Making Predictions

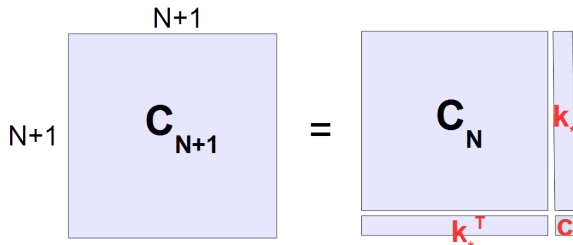
- Let's consider the joint distr. of N training responses \mathbf{y} and test response y_*

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

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



GP Regression: Making Predictions

- Given the jointly Gaussian distribution

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{C}_N & \mathbf{k}_* \\ \mathbf{k}_*^\top & c \end{bmatrix}\right)$$

GP Regression: Making Predictions

- Given the jointly Gaussian distribution

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{C}_N & \mathbf{k}_* \\ \mathbf{k}_*^\top & c \end{bmatrix}\right)$$

- 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: Making Predictions

- Given the jointly Gaussian distribution

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{C}_N & \mathbf{k}_* \\ \mathbf{k}_*^\top & c \end{bmatrix}\right)$$

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

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

- Note: The test time cost

GP Regression: Making Predictions

- Given the jointly Gaussian distribution

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{C}_N & \mathbf{k}_* \\ \mathbf{k}_*^\top & c \end{bmatrix}\right)$$

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

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

- Note: The test time cost is $\mathcal{O}(N)$

GP Regression: Making Predictions

- Given the jointly Gaussian distribution

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{C}_N & \mathbf{k}_* \\ \mathbf{k}_*^\top & c \end{bmatrix}\right)$$

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

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

- Note: The test time cost is $\mathcal{O}(N)$: linear in the number of training examples

GP Regression: Making Predictions

- Given the jointly Gaussian distribution

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{C}_N & \mathbf{k}_* \\ \mathbf{k}_*^\top & c \end{bmatrix}\right)$$

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

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

- Note: The test time cost is $\mathcal{O}(N)$: linear in the number of training examples
 - .. just like kernel SVM or nearest neighbor methods

GP Regression: Making Predictions

- Given the jointly Gaussian distribution

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{C}_N & \mathbf{k}_* \\ \mathbf{k}_*^\top & c \end{bmatrix}\right)$$

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

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

- Note: The test time cost is $\mathcal{O}(N)$: linear in the number of training examples
 - .. just like kernel SVM or nearest neighbor methods
 - However, unlike kernel SVM or nearest neighbor methods, GP also gives us the variance σ_*^2

GP Regression: Interpreting GP Predictions

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

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

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_*

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

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

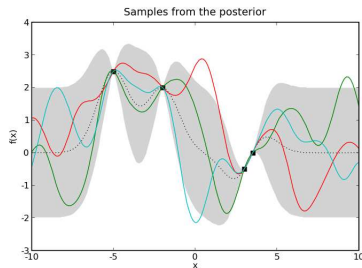
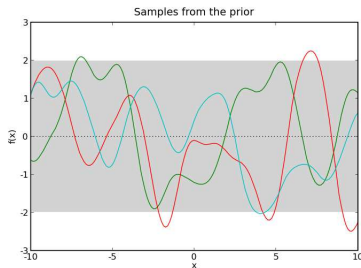
- A nearest neighbors interpretation

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

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

GP Regression: Pictorially

A GP with a squared-exponential kernel function



Left: Samples of f from the prior $\mathcal{GP}(0, \kappa)$

Right: Samples of f from the posterior of f after 4 observations

Summary

- Gaussian Process is a very powerful framework for modeling nonlinear input-output relationships
- Can think of it as Bayesian kernel methods
- Can learn the kernel hyperparameters (big advantage)
- Can use kernels that can be compositions of many kernels

Next Class

- Other examples/applications of GPs
- Inference and other issues
- Gaussian Process Latent Variable Models