

EECS 545: Machine Learning

Lecture 25. Gaussian Process Regression

Honglak Lee and Michał Dereziński

04/11/2021



Logistics

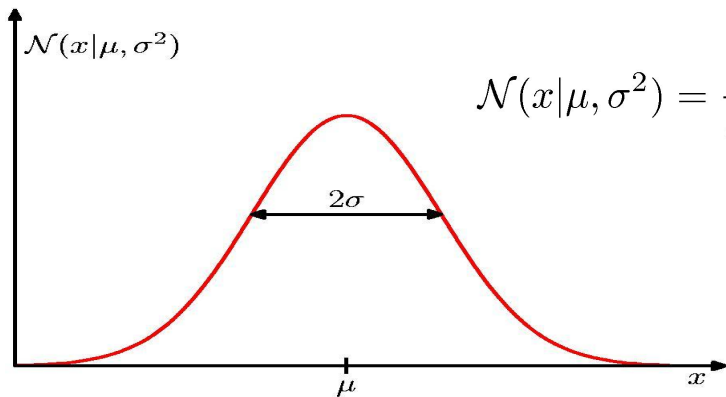
- HW6 due April 19 (last day of classes)
- Project presentations on April 21
 - In-person
 - If you are unable to be there in person, please email us and specify which group you are in
- Project final report due April 28
 - Use NeurIPS format (provided in [project info document](#), can also be [found online](#))
 - Type of content should be same / similar as those from the progress report

Outline

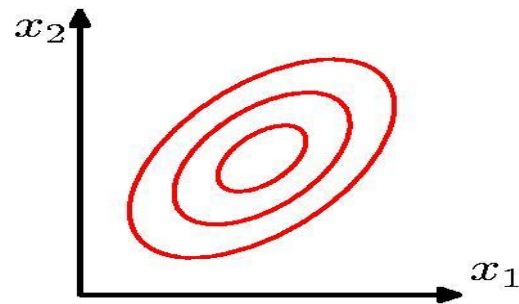
- Multivariate Gaussians
 - Marginal and conditional distributions
- Gaussian Processes
 - GP for Regression

Multivariate Gaussians – marginal and conditional distributions

The Gaussian Distribution



$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2}(x - \mu)^2 \right\}$$



$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Multivariate Gaussian (mean)

$$\begin{aligned}\mathbb{E}[\mathbf{x}] &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\} \mathbf{x} \, d\mathbf{x} \\ &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \exp \left\{ -\frac{1}{2}\mathbf{z}^T \Sigma^{-1} \mathbf{z} \right\} (\mathbf{z} + \boldsymbol{\mu}) \, d\mathbf{z} \\ &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \left(\underbrace{\exp \left\{ -\frac{1}{2}\mathbf{z}^T \Sigma^{-1} \mathbf{z} \right\} \mathbf{z}}_{\text{odd wrt } \mathbf{z}} + \underbrace{\exp \left\{ -\frac{1}{2}\mathbf{z}^T \Sigma^{-1} \mathbf{z} \right\} \boldsymbol{\mu}}_{\text{even wrt } \mathbf{z}} \right) d\mathbf{z}\end{aligned}$$

Because the integral is from $(-\infty, \infty)$,

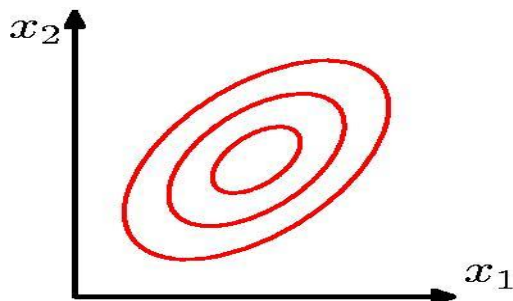
$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}$$

Multivariate Gaussian (covariance)

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma}$$

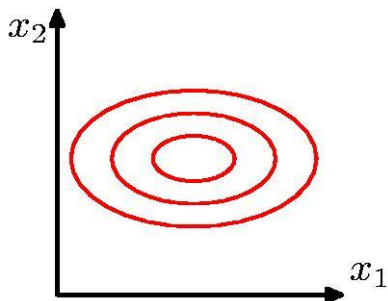
$$\text{cov}[\mathbf{x}] = \mathbb{E} [(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] = \boldsymbol{\Sigma}$$

Types of covariance matrices:



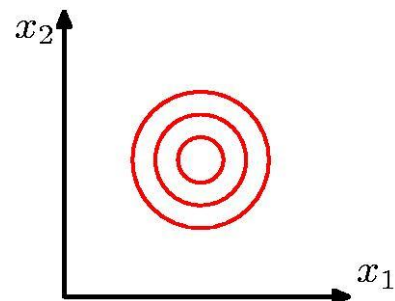
(a)

General



(b)

Diagonal



(c)

Spherical
(const*Identity)

Partitioned Gaussian Distributions

- Multivariate Gaussian distribution for \mathbf{x}

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Partitioning \mathbf{x} into \mathbf{x}_a and \mathbf{x}_b .

- Mean and covariance

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$$

- Precision matrix

$$\boldsymbol{\Lambda} \equiv \boldsymbol{\Sigma}^{-1} \quad \boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}$$

Useful Formula for Matrix Inversion

- Woodbury Matrix Inversion Lemma

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{M} & -\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{M} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \end{pmatrix}$$

- where $\mathbf{M} = (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}$.
- \mathbf{M}^{-1} is known as Schur complement of the matrix $\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}$.

Partitioned Conditionals and Marginals

- Conditional distribution

$$p(\mathbf{x}_a | \mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b})$$

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Lambda}_{aa}^{-1} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba}$$

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\Sigma}_{a|b} \{ \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \}$$

$$= \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b)$$

$$= \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} (\mathbf{x}_b - \boldsymbol{\mu}_b)$$

- Marginal distribution

$$p(\mathbf{x}_a) = \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b$$

$$= \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})$$

Derivation of conditional Gaussian

- From the exponent of the Gaussian

$$\begin{aligned} -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = \\ -\frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b). \end{aligned}$$

- Note that general Gaussians have the following form:

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1}\mathbf{x} + \mathbf{x}^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + \text{const}$$

- Treat \mathbf{x}_b as constant and rearrange terms for \mathbf{x}_a .

Derivation of conditional Gaussian

- Second order term:

$$-\frac{1}{2}\mathbf{x}_a^T \Lambda_{aa} \mathbf{x}_a \quad \Rightarrow \quad \Sigma_{a|b} = \Lambda_{aa}^{-1}.$$

- First order term:

$$\mathbf{x}_a^T \{ \Lambda_{aa} \mu_a - \Lambda_{ab} (\mathbf{x}_b - \mu_b) \}$$

– From $\Sigma_{a|b}^{-1} \mu_{a|b} = \Lambda_{aa} \mu_a - \Lambda_{ab} (x_b - \mu_b)$

$$\begin{aligned} \mu_{a|b} &= \Sigma_{a|b} \{ \Lambda_{aa} \mu_a - \Lambda_{ab} (\mathbf{x}_b - \mu_b) \} \\ &= \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} (\mathbf{x}_b - \mu_b) \end{aligned}$$

Derivation of conditional Gaussian

- Expressing in terms of covariance matrix

$$\begin{aligned}\mu_{a|b} &= \mu_a + \Sigma_{ab}\Sigma_{bb}^{-1}(\mathbf{x}_b - \mu_b) \\ \Sigma_{a|b} &= \Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba}.\end{aligned}$$

- where

$$\begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}^{-1} = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

$$\Lambda_{aa} = (\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}$$

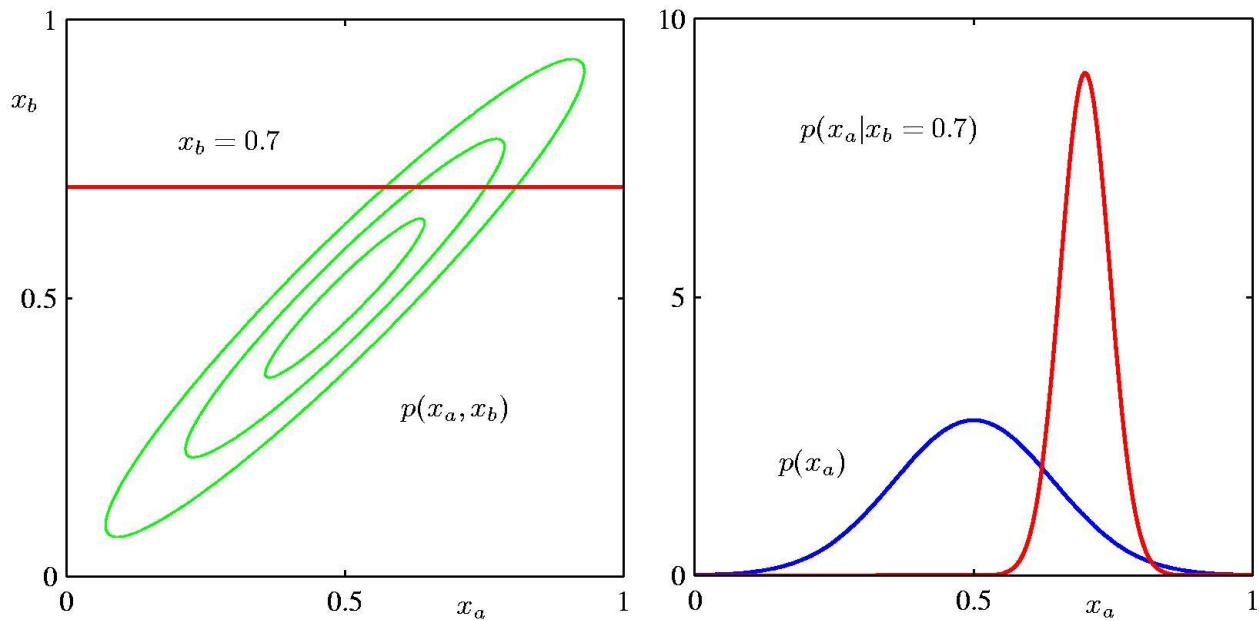
$$\Lambda_{ab} = -(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}\Sigma_{ab}\Sigma_{bb}^{-1}.$$

- Here, we used matrix inversion lemma

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{M} & -\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{M} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \end{pmatrix}$$

- where $\mathbf{M} = (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}.$

Partitioned Conditionals and Marginals



Linear Gaussian Distributions

- Linear combination of Gaussian random vector also has a Gaussian distribution

- For Gaussian random vector $X \sim \mathcal{N}(\mu_X, \Sigma_X)$, we have

$$AX + b \sim \mathcal{N}(b + A\mu_X, A\Sigma_X A^T)$$

for matrix \mathbf{A} , vector b

- Marginal and conditional distributions are also Gaussian

Bayes' Theorem for Gaussian Variables

- Given

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$$

- we arrive at the joint distribution

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{A}\boldsymbol{\mu} + \mathbf{b} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Lambda}^{-1} & \boldsymbol{\Lambda}^{-1}\mathbf{A}^T \\ \mathbf{A}(\boldsymbol{\Lambda}^{-1})^T & \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T + \mathbf{L}^{-1} \end{pmatrix}\right)$$

- and then we derive (using derivations for conditional and marginal distributions)

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T)$$

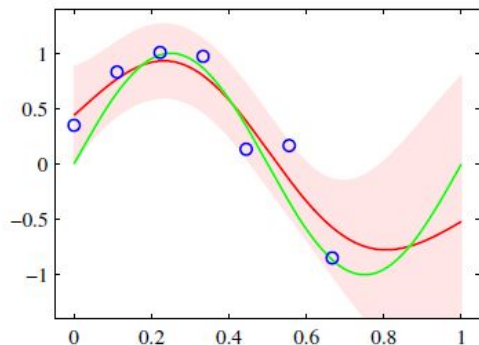
$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^T\mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma})$$

- where $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A})^{-1}$

Gaussian Processes

Why GPs?

- Here are some data points. What function did they come from?



- GPs are a nice way of expressing this “prior on functions” idea.
- Applications:
 - Regression
 - Classification

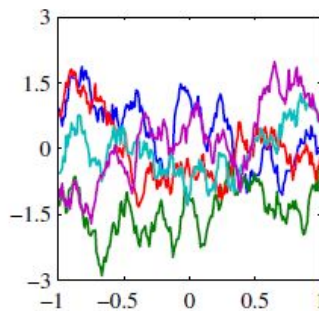
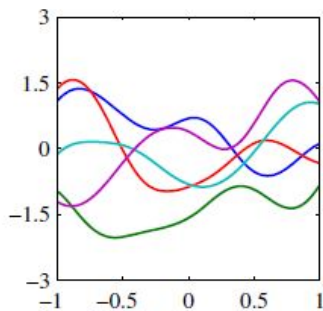
Definition of GP

- A Gaussian process is defined as a probability distribution over functions $h(\mathbf{x})$, such that the set of values of $h(\mathbf{x})$ evaluated at an arbitrary set of points $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}$ jointly have a Gaussian distribution.
 - Any finite subset of indices defines a multivariate Gaussian distribution (i.e., $h(\mathbf{x}^{(1)}), h(\mathbf{x}^{(2)}), \dots, h(\mathbf{x}^{(n)})$) is a multivariate Gaussian.
- What determines the GP is
 - The mean function $\mu(\mathbf{x}) = E(h(\mathbf{x}))$
 - The covariance function (kernel)
 $k(\mathbf{x}, \mathbf{x}') = E(h(\mathbf{x})h(\mathbf{x}'))$
 - In most applications, we take $\mu(\mathbf{x})=0$. Hence the prior is represented by the kernel.

Covariance function of GP defines prior

- The figures show samples of functions drawn from Gaussian processes for two different choices of kernel functions
 - Gaussian kernel (left) vs. exponential kernel (right)

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\theta \|\mathbf{x} - \mathbf{x}'\|_2^2)$$



$$k(\mathbf{x}, \mathbf{x}') = \exp(-\theta \|\mathbf{x} - \mathbf{x}'\|_2)$$

Linear Regression Revisited – Bayesian regression

- Bayesian linear regression model: combination of M fixed basis functions given by $\phi(\mathbf{x})$, so that

$$h(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

- Prior distribution: $p(\mathbf{w}) = N(\mathbf{w}|0, \alpha^{-1}\mathbf{I})$
- Given training data points $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}$, what is the joint distribution of $h(\mathbf{x}^{(1)}), h(\mathbf{x}^{(2)}), \dots, h(\mathbf{x}^{(n)})$?
 - \mathbf{h} is the vector with elements $h^{(i)} = h(\mathbf{x}^{(i)})$, which is given by
$$\mathbf{h} = \Phi \mathbf{w}$$
 - Where Φ is the data matrix with elements $\Phi_{nk} = \Phi_k(\mathbf{x}^{(n)})$

Linear Regression Revisited – Bayesian regression

- **$\mathbf{h} = \Phi\mathbf{w}$:** \mathbf{h} is a linear combination of Gaussian distributed variables \mathbf{w} , hence itself is Gaussian.
- Mean and covariance

$$\mathbb{E}[\mathbf{h}] = \Phi\mathbb{E}[\mathbf{w}]$$

$$\text{cov}[\mathbf{h}] = \mathbb{E}[\mathbf{h}\mathbf{h}^T] = \Phi\mathbb{E}[\mathbf{w}\mathbf{w}^T]\Phi^T = \frac{1}{\alpha}\Phi\Phi^T = K$$

where K is the Gram matrix with elements

$$K_{nm} = k(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) = \frac{1}{\alpha}\Phi(\mathbf{x}^{(n)})^T\Phi(\mathbf{x}^{(m)})$$

and $k(\mathbf{x}, \mathbf{x}')$ is the kernel function

Bayesian Linear Regression and GP

- In summary, Bayesian linear regression is a special instance of a Gaussian Process
- It is defined by the linear regression model

$$h(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

– with a weight prior

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

- The kernel function is given by

$$k(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) = \frac{1}{\alpha} \phi(\mathbf{x}^{(n)})^T \phi(\mathbf{x}^{(m)})$$

Features in Bayesian LR \Leftrightarrow kernel functions in GP.

GP for regression

- Consider the noise on the observed target values

$$y^{(n)} = h^{(n)} + \epsilon^{(n)}$$

- where $\epsilon^{(n)}$ is a random noise.

- Equivalently, consider a noise process:

$$p(y^{(n)}|h^{(n)}) = \mathcal{N}(y^{(n)}|h^{(n)}, \beta^{-1})$$

- where β is a hyperparameter (precision of the noise)

- Since $\epsilon^{(n)}$ is independent, this can be represented as multivariate Gaussian

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

GP for regression

- From the definition of GP, the marginal distribution $p(\mathbf{h})$ is given by

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

- Note: $p(\mathbf{h})$ and $p(\mathbf{y}|\mathbf{h})$ form linear Gaussian

- Then, the marginal distribution of \mathbf{y} is given by

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

- where the covariance matrix C has elements

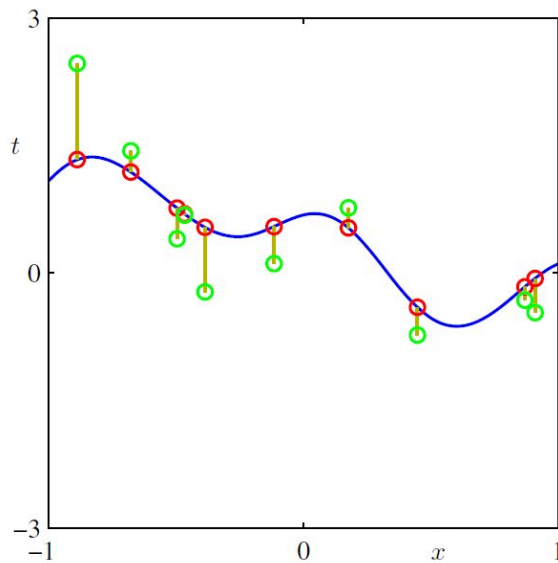
$$C(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) = k(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) + \beta^{-1}\delta_{nm}$$

Kronecker delta



Example: sampling data points

- Sample function from GP (blue): $h(\mathbf{x}) \sim GP(\mu, K)$
- Sample points from GP (red):
 $\left(\mathbf{x}^{(1)}, h(\mathbf{x}^{(1)})\right), \left(\mathbf{x}^{(2)}, h(\mathbf{x}^{(2)})\right), \left(\mathbf{x}^{(N)}, h(\mathbf{x}^{(N)})\right) \sim GP(\mu, K)$
- Add noise (green): $y^{(n)} \sim h(\mathbf{x}^{(n)}) + \mathcal{N}(0, \beta^{-1})$



GP for regression

- We have used GP to build a model of the joint distribution over sets of data points
- Goal: Given data $(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)})$ and target values $\mathbf{y} = (y^{(1)}, \dots, y^{(n)})$, predict $y^{(n+1)}$ for query point $\mathbf{x}^{(n+1)}$.
- Idea: GP assumes that

$$p(y^{(1)}, \dots, y^{(n)}, y^{(n+1)}) \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_{n+1})$$

– where \mathbf{C}_{n+1} is a $(n+1) \times (n+1)$ matrix

$$\mathbf{C}_{n+1} = \begin{pmatrix} \mathbf{C}_n & \mathbf{k} \\ \mathbf{k}^T & c \end{pmatrix}, \text{ where } \mathbf{C}_n \text{ is } n \times n \text{ matrix, and } c = k(\mathbf{x}^{(n+1)}, \mathbf{x}^{(n+1)}) + \beta^{-1}$$
$$C_n[i, j] = k(x^{(i)}, x^{(j)}) + \beta^{-1} \delta_{ij}$$

\mathbf{k} is n -dim vector where $k_i = k(x^{(n+1)}, x^{(i)})$

GP for regression

$$\begin{aligned}p(\mathbf{x}_a|\mathbf{x}_b) &= \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b}) \\ \boldsymbol{\mu}_{a|b} &= \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ \boldsymbol{\Sigma}_{a|b} &= \boldsymbol{\Lambda}_{aa}^{-1} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}\end{aligned}$$

- The conditional distribution $p(y^{(n+1)}|\mathbf{y})$ is a Gaussian distribution with mean and covariance given by

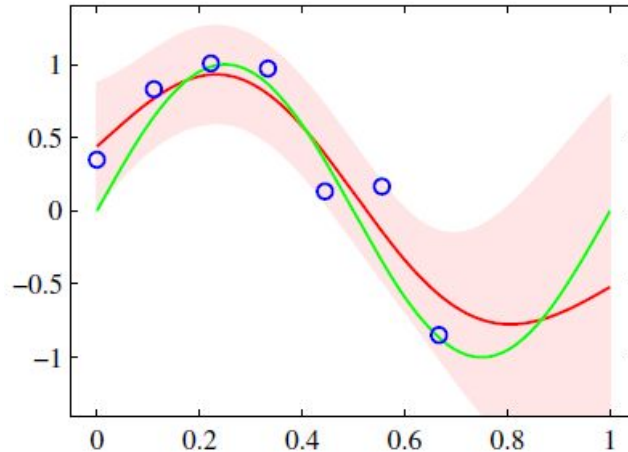
$$\begin{aligned}m(\mathbf{x}^{(n+1)}) &= \mathbf{k}^T \mathbf{C}_n^{-1} \mathbf{y} \\ \sigma^2(\mathbf{x}^{(n+1)}) &= c - \mathbf{k}^T \mathbf{C}_n^{-1} \mathbf{k}\end{aligned}$$

$$\mathbf{C}_{n+1} = \begin{pmatrix} \mathbf{C}_n & \mathbf{k} \\ \mathbf{k}^T & c \end{pmatrix}, \text{ where } \mathbf{C}_n \text{ is } n \times n \text{ matrix, and } c = k(\mathbf{x}^{(n+1)}, \mathbf{x}^{(n+1)}) + \beta^{-1}$$

- These are the key results that define Gaussian process regression.
- The predictive distribution is a Gaussian whose mean and variance both depend on $\mathbf{x}^{(n+1)}$

An Example of GP regression

- Green: underlying function (sine function)
- Blue: samples from GP (with noise)
- Red: prediction from GP regression with “error bars”



GP for Regression

- The only restriction on the kernel is that the covariance matrix given by

$$C(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) = k(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) + \beta^{-1} \delta_{nm}$$

must be positive definite.

- GP will involve a matrix of size $N \times N$, for which require $O(N^3)$ computations.

Learning Hyperparameters

- We can learn hyperparameters θ as coefficients for the kernel function
 - e.g. exponential quadratic kernel $k(\mathbf{x}^{(n)}, \mathbf{x}^{(m)}) = \theta_0 \exp \left\{ -\frac{\theta_1}{2} \|\mathbf{x}^{(n)} - \mathbf{x}^{(m)}\|^2 \right\} + \theta_2 + \theta_3 (\mathbf{x}^{(n)})^T \mathbf{x}^{(m)}$
- Log likelihood

$$\ln p(\mathbf{y}|\theta) = -\frac{1}{2} \ln |C_N| - \frac{1}{2} \mathbf{y}^T C_N^{-1} \mathbf{y} - \frac{N}{2} \ln(2\pi)$$

- Gradient Ascent for parameter θ

$$\frac{\partial}{\partial \theta_i} \ln p(\mathbf{y}|\theta) = -\frac{1}{2} \text{Tr} \left(C_N^{-1} \frac{\partial C_N}{\partial \theta_i} \right) + \frac{1}{2} \mathbf{y}^T C_N^{-1} \frac{\partial C_N}{\partial \theta_i} C_N^{-1} \mathbf{y}$$

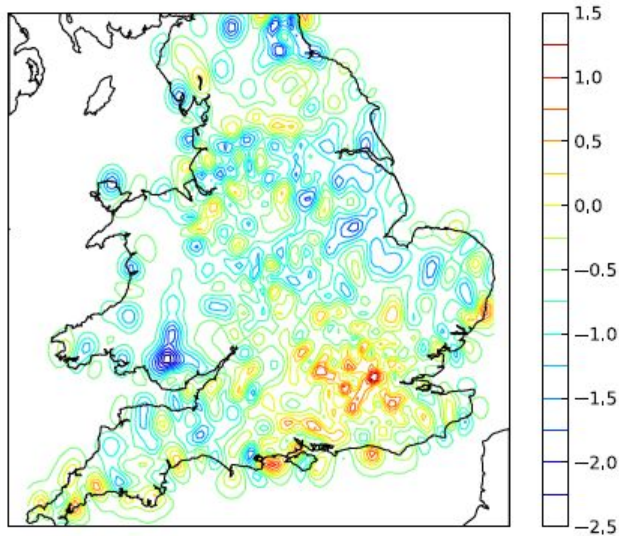
– where we used the following:

$$\frac{\partial}{\partial x} (\mathbf{A}^{-1}) = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1}$$

$$\frac{\partial}{\partial x} \ln |\mathbf{A}| = \text{Tr} \left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \right)$$

GP applications

- Modeling variability of the apartment price



James Hensman, Nicolo Fusi, and Neil D. Lawrence.

"Gaussian processes for big data." *arXiv preprint arXiv:1309.6835* (2013).

GP applications

- Modeling atmospheric CO₂ concentrations

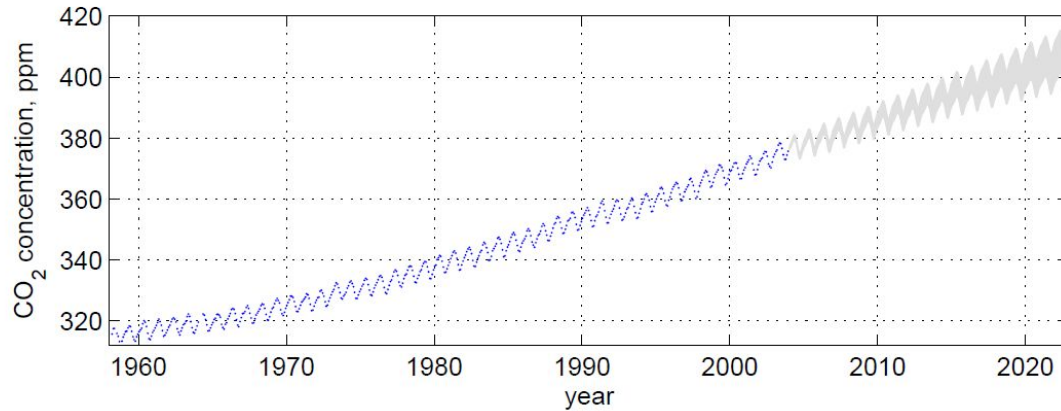


Figure 5.6: The 545 observations of monthly averages of the atmospheric concentration of CO₂ made between 1958 and the end of 2003, together with 95% predictive confidence region for a Gaussian process regression model, 20 years into the future. Rising trend and seasonal variations are clearly visible. Note also that the confidence interval gets wider the further the predictions are extrapolated.

C. E. Rasmussen & C. K. I. Williams

"Gaussian processes for machine learning." *the MIT Press* 2.3 (2006): 4.

Summary of GP

- **Practically**, you can use GP to model mean and uncertainty of your prediction.
- GP is a prior distribution over functions
- GP generates data points that are jointly a Gaussian distribution
- Most interesting structure is in $k(\mathbf{x}, \mathbf{x}')$, the kernel.
- GP can be used for regression to predict the target for a new input example.

Quiz

[Link here](#) or QR code below

