# Handout 8: Gaussian process regression

Lecturer & author: Georgios P. Karagiannis                    georgios.karagiannis@durham.ac.uk

---

**Aim.** To introduce the ideas of learning machines by introducing data into high-dimensional feature spaces for accuracy gains; introduce the kernel trick, and kernel functions.

---

**Reading list & references:**

(1) Bishop, C. M. (2006). Pattern recognition and machine learning (Vol. 4, No. 4, p. 738). New York: Springer.
- Ch. 6.4 Gaussian process
(2) Rasmussen, C. E., & Williams, C. K. (2006). Gaussian processes for machine learning (Vol. 1, p. 159). Cambridge, MA: MIT press.
- Chapter 2, Regression (supplementary)
(3) Roustant, O., Ginsbourger, D., & Deville, Y. (2012). DiceKriging, DiceOptim: Two R packages for the analysis of computer experiments by kriging-based metamodeling and optimization. Journal of statistical software, 51, 1-55.
- Supplementary material related to the the implementation of GP in R computing environment.

## 1. INTRO AND MOTIVATION

*Note* 1. As motivation for the Gaussian process regression, we "Kernelize" the standard Bayesian normal linear regression int he machine learning framework.

*Note* 2. Consider the predictive rule $h(x) = \eta(x)$, and cast it in a linear form $\eta(x) = (\psi(x))^\top w$ where $\psi(x) = (\psi_1(x), ..., \psi_d(x))$ is a vector of basis functions mapping the input space $\mathcal{X}$ into a feature space $\mathcal{F}$. Assume there is available a set of observables $\{z_i = (x_i, y_i)\}_{i=1}^n$. We associate the learning problem with the Bayesian linear regression model

(1.1)
$$
\begin{cases}
y_i | \psi(x_i), w, \sigma^2 \overset{\text{ind}}{\sim} \mathrm{N}\left(\eta(x_i), \sigma^2\right), \ i = 1, ..., n \\
\eta(\cdot) \quad\quad\quad = (\psi(\cdot))^\top w \\
w \quad\quad\quad\quad \sim \mathrm{N}(\mu_0, V_0)
\end{cases}
\text{equiv.}
\begin{cases}
y | \eta, \sigma^2 \sim \mathrm{N}\left(\eta, I\sigma^2\right) \quad \text{(sampl. distr.)} \\
\eta = \Psi w \quad\quad\quad\quad \text{(linear model restr.)} \\
w \sim \mathrm{N}(\mu_0, V_0) \quad\quad \text{(prior)}
\end{cases}
$$

where $[\Psi]_{i,j} = \psi_j(x_i)$.

*Note* 3. The marginal likelihood is

(1.2)
$$
f(y) = \mathrm{N}\left(y | \Psi^\top \mu_0, \ \Psi V_0 \Psi^\top + I\sigma^2\right)
$$

where $\mathrm{N}(y | \mu, \Sigma)$ denotes the pdf of the Normal distribution with mean $\mu$, and covariance matrix $\Sigma$.

*Note* 4. The predictive distribution of a new outcome $y_*$ at a new input $x_*$ given the observables $\{z_i = (x_i, y_i)\}_{i=1}^n$ is

$$f(y_* | x_*, \{(x_i, y_i)\}) = \mathrm{N}\left(\mu_*(x_*), \sigma_*^2(x_*)\right)$$

with

$$(1.3) \qquad \mu_*(x_*) = \psi(x_*)^\top \mu_0 + \frac{1}{\sigma^2} \overbrace{\psi(x_*)^\top V \Psi}^{K(x_*,X)=} \left(\overbrace{\Psi^\top V \Psi}^{K(X,X)=} + \sigma^2\right)^{-1} \left(\Psi^\top \mu_0 - y\right)$$

$$(1.4) \qquad \sigma_*^2(x_*) = \underbrace{\psi(x_*)^\top V \psi(x_*)}_{=K(x_*,x_*)} - \underbrace{\psi(x_*)^\top V \Psi}_{=K(x_*,X)} \left(\underbrace{\Psi^\top V \Psi}_{=K(x,X)} + \sigma^2\right)^{-1} \underbrace{\left(\psi(x_*)^\top V \Psi\right)^\top}_{=K(X,x_*)}$$

according to Proposition 38.

*Note* 5. In the prior part of (1.1), let's assume $\mu_0 = 0$ (arguably) denoting complete ignorance whether $\eta(\cdot)$ is positive or negative. By Kernel trick, in (1.3) and (1.4), the feature space always enters in the form inner products. In fact we can define a kernel $K(x, x') = \langle L\psi(x), L\psi(x')\rangle = \psi(x)^\top V \psi(x')$ where $L$ is such that $V = L^\top L$, in terms of Section 4 in Handout 7: Kernel methods. We can denote $K(x_*, X) = \psi(x_*)^\top V \Psi$, and $K(x_*, x_*) = \psi(x_*)^\top V \psi(x_*)$.

*Note* 6. No need to memorize the formulas in (1.2), (1.3), and 1.4. The material in Notes 2, 3, and 4 is given as a motivation for the Gaussian process regression.

## 2. THE GAUSSIAN PROCESS REGRESSION MODEL

**Definition 7.** Gaussian process (GP) is a collection of random variables $\{f(x); x \in \mathcal{X}\}$, indexed by label $x$, where any finite collection of those variables has a multivariate normal distribution. It is fully specified by its mean and covariance functions. It is denoted as

$$f(\cdot) \sim \mathrm{GP}(\mu(\cdot), C(\cdot, \cdot))$$

with mean

$$\mu(x) := \mathrm{E}(f(x)), \ x \in \mathcal{X}$$

and covariance function

$$C(x, x') := \mathrm{Cov}(f(x), f(x')), \ x, x' \in \mathcal{X}$$

*Note* 8. Essentially, GP is a distribution defined over functions.

*Note* 9. Consider a function $\eta : \mathcal{X} \to \mathbb{R}$ with $\eta(x) = \langle \psi(x), w \rangle$ where $\psi(x)$ is a vector of known basis (feature) functions mapping from the input space $\mathcal{X}$ to the feature space $\mathcal{F}$, and $w \in \mathbb{R}^d$ is an unknown vector a priori following a normal distribution $w \sim \mathrm{N}(0, V)$, where the prior mean is set to zero denoting complete uncertainty about the sign of $w$'s. Then the marginal $\eta(\cdot)$ follows a Normal distribution as a linear transformation of Normal variates with mean $\mathrm{E}(\eta(x)) = 0$ and covariance $\mathrm{Cov}(\eta(x), \eta(x')) = \psi(x)^\top V \psi(x')$ for any $x, x' \in \mathcal{X}$. Based on the Kernel trick and

Definition 7, we can equivalently specify $\eta\left(\cdot\right)\sim\text{GP}\left(0,C\left(\cdot,\cdot\right)\right)$ for some kernel / covariance function $C\left(x,x'\right)=\psi\left(x\right)^{\top}V\psi\left(x'\right)$.

*Note* 10. We introduce the concept of Gaussian process regression in the machine learning framework below.

*Note* 11. Consider the predictive rule $h\left(x\right)=\eta\left(x\right)$, and assume that $\eta:\mathcal{X}\rightarrow\mathbb{R}$ with unknown image (possibly up to a set of properties, we will discuss this later) and $\mathcal{X}\subseteq\mathbb{R}^d$.

**Example 12.** [1]Figure 2.1a shows a function $\eta\left(x\right)$ we pretend that we do not know but we wish to recover. To recover $\eta\left(x\right)$, we collect training data set as in Figure 2.1b.



(A) True $\eta\left(x\right)$     (B) Train-ing Examples $\{z_i=(x_i,y_i)\}$     (C) Realiza-tions of a prior Gaussian pro-cess $\eta\left(\cdot\right)$ $\sim$ $\text{GP}\left(0,C\left(\cdot,\cdot\right)\right)$
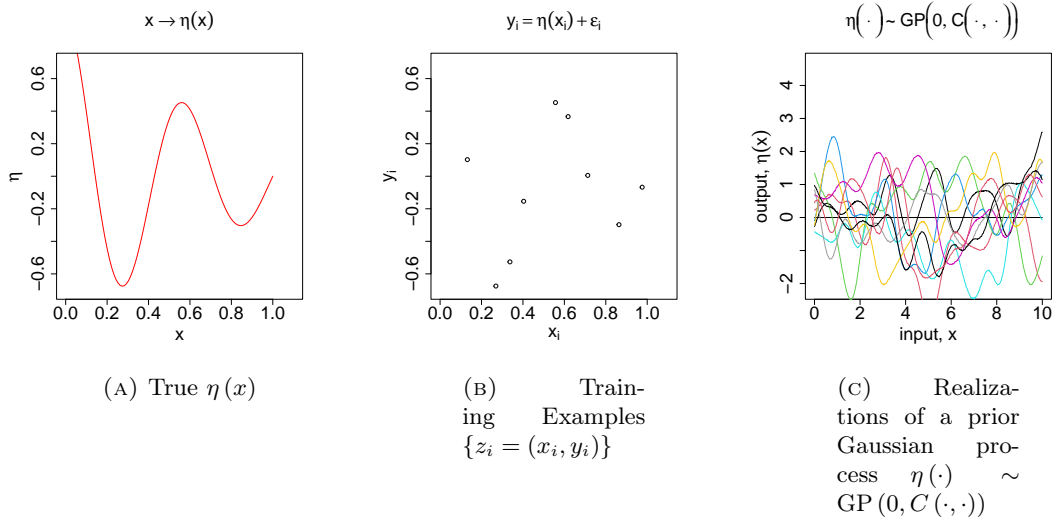
FIGURE 2.1. A toy example

*Note* 13. For training purposes, assume there is available a set of observables $\{z_i=(x_i,y_i)\}_{i=1}^n$ whose sampling distribution is such that

$$y_i=\eta\left(x_i\right)+\epsilon_i,\ \epsilon_i\overset{\text{iid}}{\sim}\text{N}\left(0,\sigma^2\right),\ i=1,...,n$$

(2.1)          or equivalently

$$y_i|\eta\left(\cdot\right),\sigma^2\overset{\text{iid}}{\sim}\text{N}\left(\eta\left(x_i\right),\sigma^2\right),\ i=1,...,n$$

for some unknown $\sigma^2>0$. (2.1) can result by considering, a quadratic loss $\ell\left(h,z=(x,y)\right)=\frac{1}{\sigma^2}\left(h\left(x\right)-y\right)^2$, and sampling distribution with pdf

$$\text{pr}\left(y|\left\{x_i,y_i\right\}\right)\propto\exp\left(-\sum_{i=1}^n\ell\left(h\left(x_i\right),\left(x_i,y_i\right)\right)\right).$$

---

As $\eta(x)$ is assumed to be unknown, according to the Bayesian paradigm and by taking advantage of Note 9, we assign a GP prior on $\eta(\cdot)$

(2.2) $$\eta(\cdot) \sim \mathrm{GP}(\mu(\cdot|\beta), C(\cdot,\cdot|\phi))$$

where $\mu$ is parametrized by unknown $\beta$ (e.g. $\mu(x|\beta) = x^\top\beta$), and $C$ is parametrized by unknown $\phi$ (e.g. $C(x,x'|\phi) = \exp\left(-\frac{1}{2\phi}\|x-x'\|_2^2\right)$; radial/Gaussian kernel). Summing up, the Bayesian model

$$\begin{cases} y_i|\eta(\cdot) & \overset{\text{iid}}{\sim} \mathrm{N}\left(\eta(x_i),\sigma^2\right), \ i = 1, ..., n \\ \eta(\cdot) & \sim \mathrm{GP}(\mu(\cdot|\beta), C(\cdot,\cdot|\phi)) \end{cases}$$

up to some unknown tuning parameters $\sigma^2$, $\beta$, and $\phi > 0$ which are suppressed from the conditioning to easy the notation.

*Note* 14. Consider $\eta_* = \eta(X_*)$ where $X_* = (x_{*,1}, x_{*,2}, ..., x_{*,m})^\top$ is a vector of new inputs of any length $m > 0$. The joint distribution of $(\eta_*, y)^\top$ is

(2.3) $$\begin{pmatrix} \eta_* \\ y \end{pmatrix} \sim \mathrm{N}\left(\begin{pmatrix} \mu(X_*) \\ \mu(X) \end{pmatrix}, \begin{pmatrix} C(X_*, X_*) + I\sigma^2 & C(X_*, X) \\ C(X, X_*) & C(X, X) + I\sigma^2 \end{pmatrix}\right)$$

where $C(X, X_*)$ is a Gram matrix over $X$ and $X_*$ such as $[C(X, X_*)]_{i,j} = C(x_i, x_{*,j})$.

*Note* 15. The conditional distribution of $\eta_* = \eta(X_*)$ given the training sample $\{z_i = (x_i, y_i)\}$, as results from 2.3 (Proposition 38),

$$\eta_*|y \sim \mathrm{N}(\mu_*(X_*), C_*(X_*, X_*))$$

is a normal distribution, with mean

(2.4) $$\mu_*(X_*) = \mathrm{E}(\eta_*|y) = \mu(X_*) + C(X_*, X)\left(C(X, X) + I\sigma^2\right)^{-1}(y - \mu(X))$$

at $X_*$ and with covariance function

(2.5) $$C_*(X_*, X_*) = \mathrm{Cov}(\eta_*|y) = C(X_*, X)\left(C(X, X) + I\sigma^2\right)^{-1}(C(X_*, X))^\top$$

*Note* 16. Because $X_*$ is of any finite length, and the derivations in Note 15, by definition of GP, the predictive distribution of $\eta(\cdot)$ given the data $\{z_i = (x_i, y_i)\}$ is the Gaussian process

(2.6) $$\eta(\cdot) \sim \mathrm{GP}(\mu_*(\cdot), C_*(\cdot,\cdot))$$

with mean function and covariance function

(2.7) $$\mu_*(x_*) = \mu(x_*) + C(x_*, X)\left(C(X, X) + I\sigma^2\right)^{-1}(y - \mu(X))$$

(2.8) $$C_*(x_*, x'_*) = C(x_*, X)\left(C(X, X) + I\sigma^2\right)^{-1} C(X, x'_*)$$

for any points $x_*, x'_* \in \mathcal{X}$. If I consider $X_* = (x_*, x'_x)^\top$, (2.7) results as the first block of 2.4, and (2.8) results as the top off-diagonal block of 2.5.
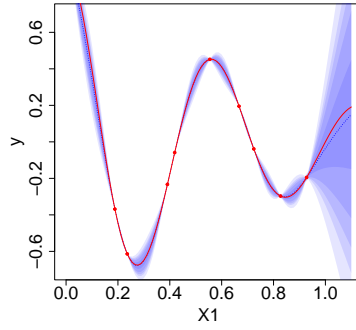
FIGURE 2.2. Predictive GP

*Note* 17. Note that the posterior expected rule at $x_* \in \mathcal{X}$ is

$$(2.9) \quad \mathrm{E}\left(h\left(x_*\right)|y\right) = \mathrm{E}\left(\eta\left(x_*\right)|y\right) = \mu\left(x_*\right) + C\left(x_*, X\right)\left(C\left(X, X\right) + I\sigma^2\right)^{-1}\left(y - \mu\left(X\right)\right)$$

$$= \sum_{i=1}^{n} \alpha_i C\left(x_i, x_*\right)$$

where $\alpha = \mu\left(x_*\right) + \left(y - \mu\left(X\right)\right)\left(C\left(X, X + I\sigma^2\right)\right)^{-1}$. This is in accordance to the Representation theorem (Theorem 21 in Handout 7 Kernel methods) with reference to the Bayesian linear regression (Note 2).

**Example 18.** Figure 2.2 presents the predictive distribution from a trained GP regression

## 3. TRAINING (VIA EMPIRICAL BAYES)

*Note* 19. Recall that the mean and covariance functions in (2.6) depend on tunable parameters $\sigma^2$, $\phi$, and $\beta$. When the number of training examples is small, the behavior of (2.6) is sensitive to these hyperparameters. In Figure 3.1, there are two instances of GP regression given 6 examples where the tunable parameters are different.
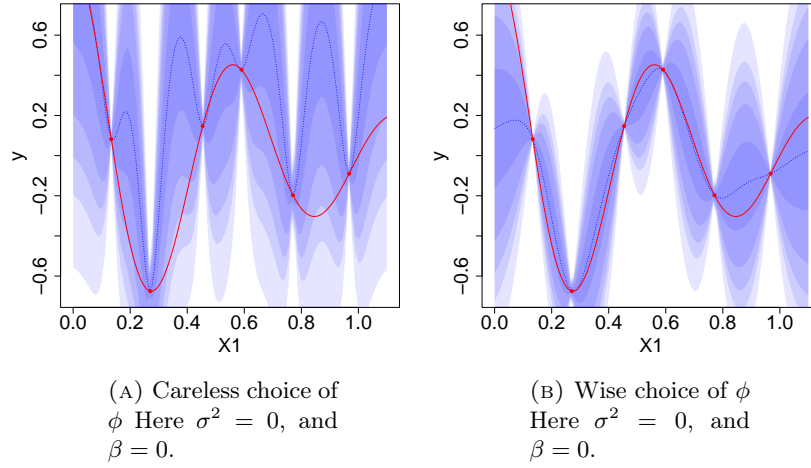
Created on 2023/03/16 at 11:36:50    by Georgios Karagiannis

(A) Careless choice of $\phi$ Here $\sigma^2 = 0$, and $\beta = 0$.

(B) Wise choice of $\phi$ Here $\sigma^2 = 0$, and $\beta = 0$.

FIGURE 3.1. Sensitivity in the tunable parameters $\phi$. Here $\sigma^2 = 0$, and $\beta = 0$.

*Note* 20. The marginal likelihood $f\left(y|\sigma^2, \phi, \beta\right)$ of $y$ given the known parameters $\sigma^2, \phi, \beta$ results from (2.3) as

$$y|\sigma^2, \phi, \beta \sim \mathrm{N}\left(\mu\left(X|\beta\right), C\left(X, X|\phi\right) + I\sigma^2\right).$$

*Note* 21. Let $\mu_\beta = \mu\left(X|\beta\right)$ and $C_\phi = C\left(X, X|\phi\right)$. To learn the unknown hyper-parameters $\theta = \left(\sigma^2, \phi, \beta\right)$ according to Empirical Bayes procedure as well as according to classical methods using as error function the marginal likelihood, we need to minimize

(3.1) $$\left(\hat{\sigma}^2, \hat{\phi}, \hat{\beta}\right) = \underset{\sigma^2, \phi, \beta}{\arg\min}\left(-2\log\left(\mathrm{N}\left(y\,|\,\mu_\beta, C_\phi + I\sigma^2\right)\right)\right)$$

$$= \underset{\sigma^2, \phi, \beta}{\arg\min}\left(\underbrace{\log\left(\left|C_\phi + I\sigma^2\right|\right) + \left(y - \mu_\beta\right)^\top\left(C_\phi + I\sigma^2\right)^{-1}\left(y - \mu_\beta\right)}_{\hat{R}(\sigma^2, \phi, \beta)}\right).$$

*Note* 22. (3.1) can be solved via GD (Algorithm 1 Handout 2 Gradient descent), or the Stochastic Gradient (Handout 3: Stochastic gradient descent) as the required gradient ca be easily computed as

No need t

memorized

$$\frac{\mathrm{d}R}{\mathrm{d}\beta_j} = \left(C_\phi + I\sigma^2\right)^{-1}\left(y - \mu_\beta\right)\frac{\mathrm{d}\mu_\beta}{\mathrm{d}\beta_j}$$

$$\frac{\mathrm{d}R}{\mathrm{d}\phi_j} = \mathrm{tr}\left(\left(C_\phi + I\sigma^2\right)^{-1}\left[\frac{\partial C_\phi}{\partial\phi_j}\right]\right) + \left(y - \mu\right)^\top\left(C_\phi + I\sigma^2\right)^{-1}\left[\frac{\partial C_\phi}{\partial\phi_j}\right]\left(C_\phi + I\sigma^2\right)^{-1}\left(y - \mu\right)$$

$$\frac{\mathrm{d}R}{\mathrm{d}\sigma^2} = \mathrm{tr}\left(\left(C_\phi + I\sigma^2\right)^{-1}\right) + \left(y - \mu\right)^\top\left(C_\phi + I\sigma^2\right)^{-1}\left(C_\phi + I\sigma^2\right)^{-1}\left(y - \mu\right)$$

4. EXAMPLES OF EXAMPLES OF COVARIANCE FUNCTIONS

*Note* 23. The covariance function $C : \mathcal{X} \times \mathcal{X} \to [0, \infty)$ with $C\left(x, x'\right)$ describes how much two random variables $x, x'$ change together.

*Note* 24. Covariance function is a functional parameter of the GP prior (2.2). Different covarince functions represent different properties, hence they impose different prior info in the GP regression model; they are crusial parameters.

*Note* 25. Any positive definite Kernel as described in Section 4 in Handout 7 Kernel methods can be used as a covarince function. Consequently kernel construction approached and theories introduced can be use for the covariance functions as well.

**Definition 26.** Stationary covariance function is called a covariance function $C : \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty)$ whose image can be written as $C(x, x') = C(\|x - x'\|)$ namely, the dependence between any pair of input points $x$ , $x'$ is a function of their distance and only.

*Note* 27. a one-dimensional Gaussian process one way to understand the characteristic length-scale of the process (if this exists) is in terms of the number of upcrossings of a level $u$. The expected number of upcrossings $\mathrm{E}(N_u)$ of the level u on the unit interval by a zero-mean, stationary, is

$$\mathrm{E}(N_u) = \frac{1}{2\pi} \sqrt{\frac{-K''(0)}{K'(0)}} \exp\left(-\frac{u^2}{2K(0)}\right)$$

*Note* 28. Popular covariance functions are

**Gaussian covariance function:** given as

$$C(r) = \exp\left(-\frac{1}{2\phi^2} r^2\right)$$

- It is infinitely differentiable, which means that the GP is very smooth.
- The parameter $\phi$ is called lengthscale.
- The number of upcrossing at level $u$ is $\mathrm{E}(N_u) = \left(2\phi^2\right)^{-1}$ meaning that smaller $\phi$ represents more upcrossings, hence represents smaller scale dependences

**Exponential Covariance Function:** given as

$$C(r) = \exp\left(-\frac{1}{\phi} |r|\right)$$

- It is not differentiable at $r = 0$

**Matern Class of Covariance Functions:** given as

$$C_\nu(r) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}r}{\phi}\right)^\nu B_\nu\left(\frac{\sqrt{2\nu}r}{\phi}\right)$$

where $B_\nu(\cdot)$ is a modified Bessel functions (description of Bessel functions is out of the scope). Matern covariance function gives the Exponential one for $\nu = 1/2$ , and the exponential one for $\nu \rightarrow \infty$.

*Note* 29. Anisotropic versions of these isotropic covariance functions can be created anisotropy by setting $r(x, x') = (x - x')^\top M (x - x')$ for some positive semi-definite matrix $M$. If $M$ is diagonal this implements the use of different length-scales on different dimension of inputs. Off-diagonal elements of $M$ implement cross-dimensional dependencies in the inputs.

*Note* 30. One may consider some low degree polynomial form $\mu\left(x|\beta\right) = \sum_{j=0}^{p} x^j \beta_j$, however in this case $\mu\left(\cdot|\beta\right)$ and $C\left(\cdot,\cdot|\phi\right)$ may compete as $C\left(\cdot,\cdot|\phi\right)$ can express such behaviors according to Note 9. For this reason, the usual specification of $\mu\left(\cdot|\beta\right)$ in (2.2) is $\mu\left(x|\beta\right) = 0$ implying a priori complete uncertainty about the sign of $\eta\left(x\right)$ at each $x$.

*Note* 31. Thinking as a statistician, one may decompose (2.2) as

$$\eta\left(\cdot\right) = \mu\left(\cdot|\beta\right) + \xi\left(\cdot|\phi\right)$$

where $\mu\left(\cdot|\beta\right)$ is modeled as a low degree polynomial (e.g., 2nd degree) representing large scale dependences (see polynomial regression in Term 1), and $\xi\left(\cdot|\phi\right) \sim \mathrm{GP}\left(0, C\left(\cdot,\cdot|\phi\right)\right)$ representing lower scale dependence by using an appropriate kernel. Seeing the big picture, (2.1) is re-states as $y_i = \mu\left(x_i|\beta\right) + \xi\left(x_i|\phi\right) + \epsilon_i$ where $\epsilon_i$ represents noise (or so short scale dependence that can be considered as noise in the model).

## 6. Practice, implementation, and code

Below is some practical examples on the implementation of Gaussian process regression in R programming environment by using the R packages: DiceKriging and DiceOptim.

- Roustant, O., Ginsbourger, D., & Deville, Y. (2012). DiceKriging, DiceOptim: Two R packages for the analysis of computer experiments by kriging-based metamodeling and optimization. Journal of statistical software, 51, 1-55.

The examples below were created for the undergraduate programme SURF 2016 at Purdue University (July 8, 2016) however they are suitable to the course.

Toy example
`https://github.com/georgios-stats/Intro_GPR_SURF_2016/blob/master/Numerical_example.ipynb`

Realistic example: The Piston Simulation function model in 2D
`https://github.com/georgios-stats/Intro_GPR_SURF_2016/blob/master/Practice_2D.ipynb`

Practice Catalytic Reaction 5D
`https://github.com/georgios-stats/Intro_GPR_SURF_2016/blob/master/Practice_CatalyticReaction_5D_solution.ipynb`

Practice Piston 7D
`https://github.com/georgios-stats/Intro_GPR_SURF_2016/blob/master/Practice_Piston_7D_solution.ipynb`

Practice Robot Arm 8D

https://github.com/georgios-stats/Intro_GPR_SURF_2016/blob/master/Practice_Robot_Arm_8D_solution.ipynbb

Created on 2023/03/16 at 11:36:50 by Georgios Karagiannis

APPENDIX A. MULTIVARIATE NORMAL DISTRIBUTION $x|\mu, \Sigma \sim \mathrm{N}_d(\mu, \Sigma)^2$

**Definition 32.** A $d$-dimensional random variable $x \in \mathbb{R}^d$ is said to have a multivariate Normal (Gaussian) distribution, if for every $d$-dimensional fixed vector $\alpha \in \mathbb{R}^d$, the random variable $\alpha^\top x$ has a univariate Normal (Gaussian) distribution.

**Definition 33.** We denote the $d$-dimensional Normal distribution with mean $\mu$ and covariance matrix $\Sigma \geq 0$ as $\mathrm{N}_d(\mu, \Sigma)$.

*Notation* 34. The $d$-dimensional standardized Normal distribution is $\mathrm{N}_d(0, I)$.

**Proposition 35.** *Let random variable $x \sim N_d(\mu, \Sigma)$ , fixed vector $c \in \mathbb{R}^q$ and fixed matrix $A \in \mathbb{R}^q \times \mathbb{R}^d$. The random vector $y = c + Ax$ has distribution $y \sim N_q\left(c + A\mu, A\Sigma A^\top\right)$.*

**Proposition 36.** *Let a $d$-dimensional random vector $x \sim N_{(any)}(\mu, \Sigma)$.*
  (1) *Let $y = Ax$ and $z = Bx$, where $A \in \mathbb{R}^{q \times d}$ and $B \in \mathbb{R}^{k \times d}$: The vectors $y = Ax$ and $z = Bx$ are independent if and only if $A\Sigma B^\top = 0$.*
  (2) *Let $x = (x_1, ..., x_d)^\top$: The $x_1, ..., x_d$ are mutually independent if and only if the corresponding off diagonal parts of the $\Sigma$ are zero.*

**Proposition 37.** *Any sub-vector of a vector with multivariate Normal distribution has a multivariate Normal distribution.*

**Proposition 38.** *[Marginalization & conditioning] Let $x \sim \mathrm{N}_d(\mu, \Sigma)$. Consider partition such that*

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} ; \qquad \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} ; \qquad \Sigma = \begin{bmatrix} \Sigma_1 & \Sigma_{21}^\top \\ \Sigma_{21} & \Sigma_2 \end{bmatrix},$$

*where $x_1 \in \mathbb{R}^{d_1}$, and $x_2 \in \mathbb{R}^{d_2}$ Then:*
  (1) *For the marginal, it is $x_1 \sim \mathrm{N}_{d_1}(\mu_1, \Sigma_1)$.*
  (2) *For the conditional, if $\Sigma_1 > 0$, it is*

$$x_2|x_1 \sim \mathrm{N}_{d_2}(\mu_{2|1}, \Sigma_{2|1})$$

*where*

(A.1) $$\mu_{2|1} = \mu_2 + \Sigma_{21}\Sigma_1^{-1}(x_1 - \mu_1) \ \text{ and } \ \Sigma_{2|1} = \Sigma_2 - \Sigma_{21}\Sigma_1^{-1}\Sigma_{21}^\top$$

**Proposition 39.** *The density function of the $d$-dimensional Normal distribution with mean $\mu$ and covariance matrix $\Sigma$, when $\underline{\Sigma}$ is symmetric positive definite matrix $(\Sigma > 0)$, exists and it is equal to*

(A.2) $$f(x) = (2\pi)^{-\frac{d}{2}} \det(\Sigma)^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right)$$

---

[2] More detailed material about the Multivariate Normal distribution can be found in the can be found in "Handout 2: Revision in mixture of probability distributions" of the module "Bayesian Statistics III/IV (MATH3341/4031)" Michaelmas term, 2021 available from `https://github.com/georgios-stats/Bayesian_Statistics_Michaelmas_2021/blob/main/Lecture_handouts/02_Revision_in_mixture_of_probability_distributions.pdf`. The material in this section is just a sub-set of the statements in the referenced handout.

Page 10                    Created on 2023/03/16 at 11:36:50                    by Georgios Karagiannis