

# Machine Learning

## Gaussian Processes

---

Wojciech Czech

January 15, 2024

Institute of Computer Science

# Table of contents

1. Learning hyperparameters
2. Gaussian process for classification
3. Laplace approximation
4. Laplace approximation for Gaussian process classification
5. Automatic relevance determination
6. Gaussian processes - summary

# Learning hyperparameters

---

# Fitting Gaussian process kernel

- Create parametrized kernel function by addition or multiplication of basic kernels to reflect different characteristics of the data:
  - a. Long term smooth change
  - b. Seasonality
  - c. Short term irregularities
  - d. Medium term irregularities
  - e. Precision of white noise
  - f. Length scale of correlations
- Learn hyperparameters of a kernel function by maximizing log likelihood function

## Parametrized kernel functions

$$k(\mathbf{x}_i, \mathbf{x}_j) = \theta_0 \exp \left( -\frac{\theta_1}{2} \|\mathbf{x}_i - \mathbf{x}_j\|^2 \right) + \theta_2 + \theta_3 \mathbf{x}_i^T \mathbf{x}_j \quad (1)$$

$$\begin{aligned} k(\mathbf{x}_i, \mathbf{x}_j) = & \theta_0 \exp \left( -\frac{\theta_1}{2} \|\mathbf{x}_i - \mathbf{x}_j\|^2 \right) + \\ & \theta_2 \exp \left( -\frac{2}{\theta_3} \sin^2 \left( \pi \frac{\|\mathbf{x}_a - \mathbf{x}_b\|}{\theta_4} \right) \right) \exp \left( -\frac{\theta_5}{2} \|\mathbf{x}_a - \mathbf{x}_b\|^2 \right) + \\ & \theta_6 \left( 1 + \frac{\|\mathbf{x}_a - \mathbf{x}_b\|^2}{2\theta_7\theta_8} \right)^{-\theta_7} \end{aligned} \quad (2)$$

# Log Likelihood function

$$t(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \quad (3)$$

$$p(\mathbf{t}|\mathbf{X}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{t}|\boldsymbol{\mu}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}}) = \mathcal{N}(\mathbf{t}|m(\mathbf{X}), \mathbf{C}_{\boldsymbol{\theta}}) \quad (4)$$

$$m(\mathbf{X}) = \mathbf{0} \quad (5)$$

$$\boldsymbol{\theta}_{ML} = \operatorname{argmax}_{\boldsymbol{\theta}} \{p(\mathbf{t}|\mathbf{X}, \boldsymbol{\theta})\} = \operatorname{argmin}_{\boldsymbol{\theta}} \{-\ln p(\mathbf{t}|\mathbf{X}, \boldsymbol{\theta})\} \quad (6)$$

$$\ln p(\mathbf{t}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2} \ln |\mathbf{C}_{\boldsymbol{\theta}}| - \frac{1}{2} \mathbf{t}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{t} - \frac{n}{2} \ln(2\pi) \quad (7)$$

## Maximum likelihood solution

$$\frac{\partial}{\partial \theta_i} \ln p(\mathbf{t}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2} \text{Tr} \left( \mathbf{C}_{\boldsymbol{\theta}}^{-1} \frac{\partial \mathbf{C}_{\boldsymbol{\theta}}}{\partial \theta_i} \right) + \frac{1}{2} \mathbf{t}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \frac{\partial \mathbf{C}_{\boldsymbol{\theta}}}{\partial \theta_i} \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{t} \quad (8)$$

$$\frac{\partial}{\partial \theta_i} \ln p(\mathbf{t}|\mathbf{X}, \boldsymbol{\theta}) = 0 \quad (9)$$

$\ln p(\mathbf{t}|\mathbf{X}, \boldsymbol{\theta})$  is non-convex and can have multiple local maxima

## Practical approach

Fit the parameters  $\theta$  by maximizing likelihood  $p(\mathbf{t}|\mathbf{X}, \theta)$  of the Gaussian process distribution based on observed data  $(\mathbf{X}, \mathbf{t})$ .

1. Sum all kernels to single parametrized kernel containing all signal characteristics
2. Divide observed data into mini-batches
3. Define log likelihood function
4. Use *Adam Optimizer*[3] to minimize negative log likelihood (run optimization on all batches)
5. Use fitted kernel, mean function and observed data to make posterior predictions on unobserved data



# Adam optimizer

- Stochastic gradient descent method based on adaptive estimation of first-order and second-order moments (**AD**Aptive **M**oment estimation)
- `tf.keras.optimizers.Adam`
- State-of-the-art optimization method well-suited to large parameter spaces
- Adam combines the best properties of the AdaGrad and RMSProp algorithms to provide an optimization algorithm that can handle sparse gradients on noisy problems [2]
- Adam is relatively easy to configure where the default configuration parameters do well on most problems [2]

# Gaussian process for classification

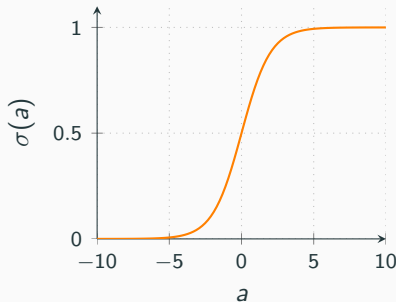
---

# Probabilistic two-class classification problem

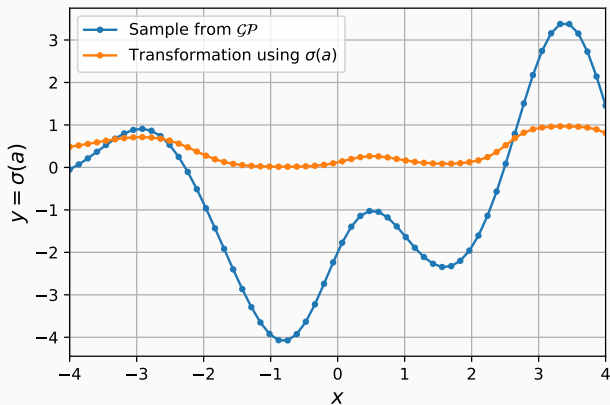
- Target:  $t \in \{0, 1\}$ , which follows Bernoulli distribution
- Define Gaussian process over function  $a(\mathbf{x})$ , then transform function using logistic sigmoid  $y = \sigma(a)$ , where:

$$\sigma(a) = \frac{1}{1 + \exp(-a)} \quad (10)$$

- We obtain non-Gaussian stochastic process over functions  $y(\mathbf{x})$ ,  $y \in (0, 1)$



# Logistic sigmoid transformation



$$p(t|a) = \sigma(a)^t (1 - \sigma(a))^{(1-t)} \quad (11)$$

# Classification problem

- Training data
  - a) Target values:  $\mathbf{t}_n = [t_1, \dots, t_n]^T$ ,  $t_i \in \{0, 1\}$
  - b) Input values:  $\mathbf{X}_n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$
- Test data:  $\mathbf{x}_{n+1}$  with target variable  $t_{n+1}$
- Determine predictive distribution  $p(t_{n+1} | \mathbf{t}_n, \mathbf{X}_n, \mathbf{x}_{n+1})$
- Gaussian process prior over vector  $\mathbf{a}_{n+1} = [a(\mathbf{x}_1), \dots, a(\mathbf{x}_{n+1})]$

$$p(\mathbf{a}_{n+1}) = \mathcal{N}(\mathbf{a}_{n+1} | \mathbf{0}, \mathbf{C}_{n+1}) \quad (12)$$

$\mathbf{C}_{n+1} \in \mathbb{R}^{(n+1) \times (n+1)}$  - covariance matrix

$$[\mathbf{C}_{n+1}]_{i,j} = C(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j) + \nu \delta_{ij} \quad (13)$$

# Classification problem - predictive distribution

- Kernel function depends on vector of parameters  $\theta$
- Non-Gaussian stochastic process over  $\mathbf{t}_{n+1}$
- Conditioning on training data (predictive distribution):

$$p(t_{n+1} = 1 | \mathbf{t}_n) = \int p(t_{n+1} = 1 | a_{n+1}) p(a_{n+1} | \mathbf{t}_n) da_{n+1} \quad (14)$$

$$p(t_{n+1} = 1 | a_{n+1}) = \sigma(a_{n+1}) \quad (15)$$

- Predictive distribution cannot be derived analytically

# Laplace approximation

---

# Laplace approximation

## Laplace approximation

Gaussian approximation to a probability density defined over a set of continuous variables (see [1]).

$$p(z) = \frac{1}{Z} f(z) \quad (16)$$

$$Z = \int f(z) dz \quad (17)$$

## Single continuous variable

Find Gaussian approximation  $q(z)$ , which is centered on a mode of the distribution  $p(z)$ .



## Mode of $p(z)$

- Find a point  $z_0$ , so that  $f'(z_0) = 0$
- Taylor expansion of  $u(z) = \ln f(z)$

$$u(z) \simeq u(z_0) + \underbrace{u'(z_0)(z - z_0)}_0 + u''(z_0)(z - z_0)^2 \quad (18)$$

$$f(z) \simeq f(z_0) \exp \left\{ -\frac{A}{2}(z - z_0)^2 \right\} \quad (19)$$

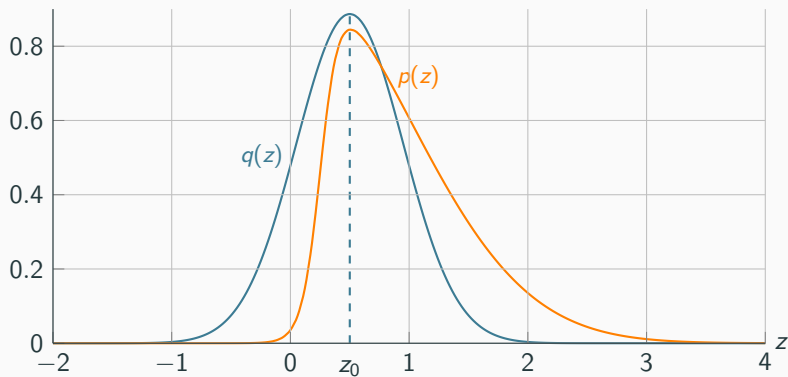
$$A = -u''(z_0) \quad (20)$$

- Normalized distribution  $q(z)$

$$q(z) = \left( \frac{A}{2\pi} \right)^{1/2} \exp \left\{ -\frac{A}{2}(z - z_0)^2 \right\} \quad (21)$$

- Gaussian approximation will only be well defined if  $A > 0$   
(maximum)

## Example 1: Visualization of the Laplace approximation



# Multidimensional Laplace approximation

$$p(\mathbf{z}) = \frac{1}{Z} f(\mathbf{z}) \quad (22)$$

$$Z = \int f(\mathbf{z}) d\mathbf{z} \quad (23)$$

- Find a point  $\mathbf{z}_0$ , so that  $\nabla f(\mathbf{z}_0) = \mathbf{0}$
- Taylor expansion of  $u(\mathbf{z}) = \ln f(\mathbf{z})$

$$u(\mathbf{z}) \simeq u(\mathbf{z}_0) - \frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \quad (24)$$

- Hessian matrix  $\mathbf{A} \in \mathbb{R}^{m \times m}$

$$\mathbf{A} = -\nabla^2 u(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0} = -\Delta \ln f(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0} \quad (25)$$

# Multidimensional Laplace approximation - taking exponential

$$f(\mathbf{z}) \simeq f(\mathbf{z}_0) \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\} \quad (26)$$

- Normalized distribution  $q(\mathbf{z})$

$$q(\mathbf{z}) = \frac{|\mathbf{A}|^{1/2}}{(2\pi)^{m/2}} \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\} \quad (27)$$

- Gaussian approximation will be only well defined if  $\mathbf{A}$  is positive definite (stationary point  $\mathbf{z}_0$  must be local maximum)

# Laplace approximation - summary

1. Find mode  $\mathbf{z}_0$  (optimization algorithm)
2. Evaluate Hessian matrix at  $\mathbf{z}_0$
3. In case of multi-modal distribution  $p(\mathbf{z})$  obtain multiple Laplace approximations
4. Laplace approximation works well when the number of data points is large (central limit theorem)
5. Approximation is local - distribution at specific value of variable

# **Laplace approximation for Gaussian process classification**

---

$$p(t_{n+1} = 1 | \mathbf{t}_n) = \int \sigma(a_{n+1}) p(a_{n+1} | \mathbf{t}_n) da_{n+1} \quad (28)$$

- Approximation of the convolution of a logistic sigmoid with a Gaussian distribution [1]:

$$\int \sigma(a) \mathcal{N}(a | \mu, \sigma^2) da \simeq \sigma(\kappa(\sigma^2)\mu) \quad (29)$$

$$\kappa(\sigma^2) = (1 + \pi\sigma^2/8)^{-1/2} \quad (30)$$

- Gaussian approximation of posterior distribution  $p(a_{n+1} | \mathbf{t}_n)$  needed

## Posterior distribution and Bayes' theorem

$$\begin{aligned} p(a_{n+1}|\mathbf{t}_n) &= \int p(a_{n+1}, \mathbf{a}_n|\mathbf{t}_n) d\mathbf{a}_n \\ &= \frac{1}{p(\mathbf{t}_n)} \int p(a_{n+1}, \mathbf{a}_n) p(\mathbf{t}_n|a_{n+1}, \mathbf{a}_n) d\mathbf{a}_n \\ &= \frac{1}{p(\mathbf{t}_n)} \int p(a_{n+1}|\mathbf{a}_n) p(\mathbf{a}_n) p(\mathbf{t}_n|\mathbf{a}_n) d\mathbf{a}_n \\ &= \int \underbrace{p(a_{n+1}|\mathbf{a}_n)}_{\mathcal{GP} \text{ regression}} \underbrace{p(\mathbf{a}_n|\mathbf{t}_n)}_{\Delta \text{ approximation}} d\mathbf{a}_n \end{aligned} \tag{31}$$



$$p(a_{n+1}|\mathbf{a}_n) = \mathcal{N}(a_{n+1} | \mathbf{k}^T \mathbf{C}_n^{-1} \mathbf{a}_n, c - \mathbf{k}^T \mathbf{C}_n^{-1} \mathbf{k}_n) \quad (32)$$

## Laplace approximation of posterior distribution $p(\mathbf{a}_n|\mathbf{t}_n)$

- Prior  $p(\mathbf{a}_n)$  is given by zero-mean Gaussian process with covariance matrix  $\mathbf{C}_n$
- Assumption about independence of training data points gives:

$$p(\mathbf{a}_n|\mathbf{t}_n) = \prod_{i=1}^n \sigma(a_i)^{t_i} (1 - \sigma(a_i))^{1-t_i} = \prod_{i=1}^n e^{a_i t_i} \sigma(-a_i) \quad (33)$$

- Taylor expansion of  $\ln p(\mathbf{a}_n|\mathbf{t}_n)$ :

$$\begin{aligned} \Psi(\mathbf{a}_n) &= \ln p(\mathbf{a}_n) + \ln p(\mathbf{a}_n|\mathbf{t}_n) \\ &\simeq -\frac{1}{2} \mathbf{a}_n^T \mathbf{C}_n^{-1} \mathbf{a}_n - \frac{n}{2} \ln(2\pi) - \frac{1}{2} |\mathbf{C}_n| + \mathbf{t}_n^T \mathbf{a}_n \\ &\quad - \sum_{i=1}^n \ln(1 + e^{a_i}) + \text{const.} \end{aligned} \quad (34)$$

## Laplace approximation of posterior distribution $p(\mathbf{a}_n | \mathbf{t}_n)$

- Find mode of posterior distribution  $p(\mathbf{a}_n | \mathbf{t}_n)$
- Gradient of  $\Psi(\mathbf{a}_n)$ :

$$\nabla \Psi(\mathbf{a}_n) = \mathbf{t}_n - \boldsymbol{\sigma}_n - \mathbf{C}_n^{-1} \mathbf{a}_n \quad (35)$$

$$\boldsymbol{\sigma} = [\sigma(a(\mathbf{x}_1)), \dots, \sigma(a(\mathbf{x}_n))] = [\sigma(a_1), \dots, \sigma(a_n)]$$

- Laplacian of  $\Psi(\mathbf{a}_n)$ :

$$\nabla^2 \Psi(\mathbf{a}_n) = -\mathbf{W}_n - \mathbf{C}_n^{-1} \quad (36)$$

$$\mathbf{W}_n = \text{diag}(\sigma(a_i)(1 - \sigma(a_i)))$$

- Hessian matrix  $\mathbf{A} = -\nabla^2 \Psi(\mathbf{a}_n)$  is positive definite - posterior distribution has a single global maximum

$$\mathbf{a}_n^* = \mathbf{C}_n(\mathbf{t}_n - \boldsymbol{\sigma}_n) \quad (37)$$

- Gaussian approximation to the posterior  $p(\mathbf{a}_n|\mathbf{t}_n)$  is given by:

$$q(\mathbf{a}_n) = \mathcal{N}(\mathbf{a}_n|\mathbf{a}_n^*, \mathbf{A}^{-1}) \quad (38)$$

# Gaussian process classification - summary

**Goal:** estimate predictive distribution:

$$p(t_{n+1} = 1 | \mathbf{t}_n) = \int \sigma(a_{n+1}) p(a_{n+1} | \mathbf{t}_n) da_{n+1} \quad (39)$$

1. Resolve Gaussian process regression:

$$p(a_{n+1} | \mathbf{a}_n) = \mathcal{N}(a_{n+1} | \mathbf{k}^T \mathbf{C}_n^{-1} \mathbf{a}_n, c - \mathbf{k}^T \mathbf{C}_n^{-1} \mathbf{k}_n) \quad (40)$$

2. Find Laplace approximation of posterior distribution  $p(\mathbf{a}_n | \mathbf{t}_n)$

$$q(\mathbf{a}_n) = \mathcal{N}(\mathbf{a}_n | \mathbf{a}_n^*, \mathbf{A}^{-1}) \quad (41)$$

3. Use Bayes' theorem (Eq. 31) and equations Eq. 32, Eq. 38 to calculate:

$$p(a_{n+1} | \mathbf{t}_n) = \int p(a_{n+1} | \mathbf{a}_n) p(\mathbf{a}_n | \mathbf{t}_n) d\mathbf{a}_n \quad (42)$$

4. Use equations Eq. 29 and Eq. 30 to calculate convolution given by Eq. 39 and get predictive distribution

# Gaussian process classification - kernel parameters

- Determine kernel function parameters  $\theta$  using maximum likelihood approach for  $p(\mathbf{t}_n|\theta)$

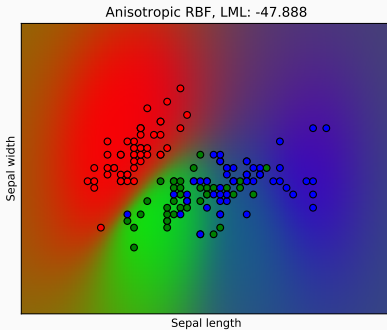
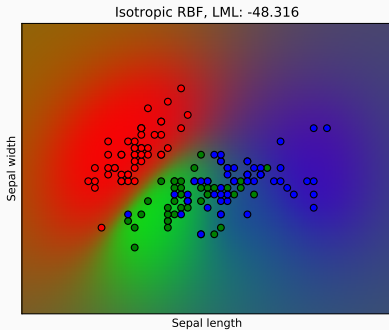
$$p(\mathbf{t}_n|\theta) = \int p(\mathbf{t}_n|\mathbf{a}_n)p(\mathbf{a}_n|\theta)d\mathbf{a}_n \quad (43)$$

- Laplacian approximation can be used again
- Approximation of log likelihood function

$$\ln p(\mathbf{t}_n|\theta) = \Psi(\mathbf{a}_n^*) - \frac{1}{2} \ln |\mathbf{W}_n + \mathbf{C}_n^{-1}| + \frac{n}{2} \ln 2\pi \quad (44)$$

- Use standard nonlinear optimization algorithms to determine a value for  $\theta$

## Example 2: Iris classification using Gaussian process



## Example 2: Configuration

- Gaussian process classification (GPC) based on Laplace approximation (Scikit Learn)
- Two types of kernels tested: Isotropic RBF and Anisotropic RBF
- Optimal kernel parameters  $\theta$  are selected during fitting
- Maximal log-likelihood value reported
- Probabilities represented using RGB heatmaps



$$k(\mathbf{x}, \mathbf{x}') = \exp \left( -(\mathbf{x} - \mathbf{x}')^T \mathbf{M} (\mathbf{x} - \mathbf{x}') \right) \quad (45)$$

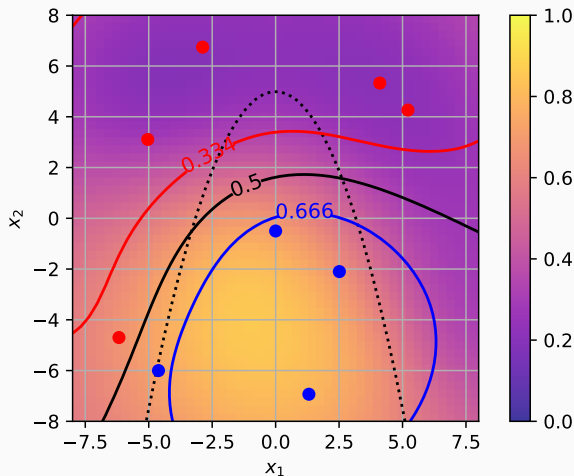
$$\mathbf{M} = \text{diag}(\boldsymbol{\theta}) \quad (46)$$

$$\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{m \times 1}$$

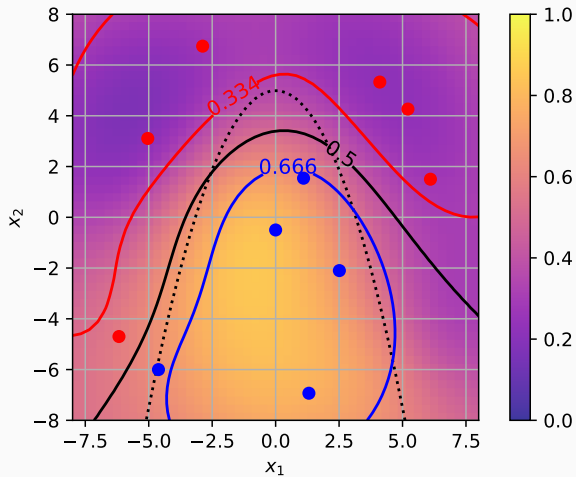
$$\boldsymbol{\theta} \in \mathbb{R}^{m \times 1}$$

$$\mathbf{M} \in \mathbb{R}^{m \times m}$$

### Example 3: RBF kernel iso-probability lines



### Example 3: RBF kernel iso-probability lines



# **Automatic relevance determination**

---

- Construct special kernel function with parameter for each input variable
- Optimize parameters on training set using maximum likelihood
- Get information about relative importance of each dimension in training set
- Detect input variables that have little effect on the predictive distribution (discard them)

$$k(\mathbf{x}_i, \mathbf{x}_j) = \theta_0 \exp \left\{ -\frac{1}{2} \sum_{k=1}^D \eta_k (x_{ik} - x_{jk})^2 \right\} \quad (47)$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \theta_0 \exp \left\{ -\frac{1}{2} \sum_{k=1}^D \eta_k (x_{ik} - x_{jk})^2 \right\} + \theta_1 + \theta_2 \sum_{k=1}^D x_{ik} x_{jk} \quad (48)$$

## Gaussian processes - summary

---

# Advantages

- Parameter tuning within the framework (easier to select hyperparameters than in case of linear regression or neural networks)
- Predictions are probabilistic (confidence intervals)
- Flexibility in selecting prior (kernel)
- Essentially non-parametric model (black box)
- Convenient probabilistic framework well suited for many tasks



# Disadvantages

- Expensive computationally (optimization of parameter space)
- $\mathcal{O}(n^3)$
- Gaussian processes are not sparse - they use the whole samples/features information to build prediction model
- Gaussian processes lose efficiency in high dimensional spaces ( $> 10^5$ )



C. M. Bishop.

***Pattern Recognition and Machine Learning.***

Springer, 2006.



J. Brownlee.

***Gentle Introduction to the Adam Optimization Algorithm for Deep Learning, Accessed January, 2022.***

Available at <https://machinelearningmastery.com/adam-optimization-algorithm-for-deep-learning/>.



D. P. Kingma and J. Ba.

**Adam: A method for stochastic optimization.**

*arXiv preprint arXiv:1412.6980*, 2014.