# 3CS Algorithm for Efficient Gaussian Process Model Retrieval

**Fabian Berns[1], Kjeld Schmidt[1], Ingolf Bracht[1], Christian Beecks[1,2]**

[1]University of Münster
Department of Computer Science -- Data Management and Analytics Group

[2]Fraunhofer Institute for Applied Information Technology FIT

living.knowledge

# Introduction (I)

- Gaussian Process Models (GPM) have been widely applied for various pattern recognition tasks

- We identified three kinds of GPMs

    i. **Default instantiation**

    ii. **Domain-specific instantiations**

    iii. **Automatically retrieved GPM given arbitrary data**

- GPM Evaluation and Application usually suffers from $\mathcal{O}(n^3)$ complexity

# Automatic GPM Retrieval – Formalization

- Gaussian Process: $f \sim GP(m, k)$

- Marginalizing $GP_\theta$ for certain data $D$ yields:
$$\mathcal{L}(m, k, \theta | D) = \frac{1}{2} \cdot [(y - \mu)^T K^{-1}(y - \mu) + \log|K| + n2\pi],$$

- Goal: find best $GP_\theta^* \in \mathbb{G}$ for $D$
$$\mathbb{G} = \{GP_\theta(m, k) | m \in \mathbb{R}^{\mathcal{X}}, k \in \mathbb{R}^{\mathcal{X} \times \mathcal{X}} \in \mathbb{R}, \theta \in \Theta\}$$

- Automatic GPM Retrieval is defined as follows:
$$GP_\theta^* = \mathrm{argmax}_{GP_\theta^*(m,k) \in \mathcal{G}} \mathcal{L}(m, k, \theta | D), \mathcal{G} \subset \mathbb{G}$$

# Concatenated Composite Covariance Search (3CS)

Covariance function is partitioned by means of change points
$T = \{\tau_i\}_{i=1}^a$:

$$\mathcal{K}(x, x'|\{k_i\}_{i=1}^a, T) = \sum_{i=1}^a k_i(x, x') \cdot 1_{\tau_{i-1} < x \leq \tau_i}(x) \cdot 1_{\tau_{i-1} < x' \leq \tau_i}(x')$$

We define the virtual search space as follows:

$$\mathcal{G}_T = \{GP_\theta(m, k)|m \in 0^{\mathcal{X}}, k = \mathcal{K}(\cdot, \cdot \,|\{k_i|k_i \in \mathcal{S}\}_{i=1}^a, T), \theta \in \Theta\}$$

$$\mathcal{S} = \left\{ \sum \prod b \,|b \in \mathcal{B} \right\}$$

**Algorithm 1** 3CS

1: **function** $(D, \mathcal{B}, c, w)$
2:     $K = \emptyset, T = \emptyset$
3:     left $= 0$, right $= max(w, n)$
4:     **while** left $< n$ **do**
5:         $D_i = \{X[\text{left, right}], Y[\text{left, right}]\}$
6:         $\tau^* = \arg\max_{\tau \in X_i} \mathcal{L}(\mathcal{K}(\cdot, \cdot \,|\{k_{\text{WN}}^l, k_{\text{WN}}^r\}, \{\tau\}) \,|\, D_i)$
7:         **if** $\tau^* \neq$ left $\wedge \tau^* \neq$ right **then**
8:             $D_i = \{X[\text{left}, \tau^*], Y[\text{left}, \tau^*]\}$
9:             $k^* = \arg\max_{k \in \mathcal{C}} \mathcal{L}(k \mid D_i)$
10:        $K = K \cup \{k^*\}, T = T \cup \{\tau^*\}$
11:        left $= \tau^*$, right $= \tau^* + w$
12:     **else**
13:        right $=$ right $+ w$
14:     **end if**
15:     **end while**
16:     $T = T \cup \{x_1, x_n\}$
17:     **return** $\mathcal{K}(\cdot, \cdot \,|\, K, T)$
18: **end function**

# Evaluation

- Eight Benchmark Datasets were used

  - **144 – 2M data records**

- 3CS **outperforms state-of-the-art algorithms** with regards to runtime and model accuracy

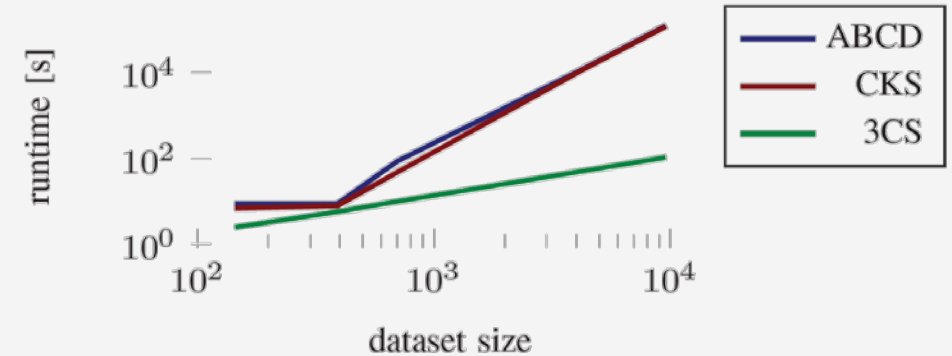- 3CS proves to be **scalable** to large datasets, while maintaining model quality



Fig. 4. Runtime for different state-of-the-art algorithms.

| Dataset | Runtime | | MSE |
|---|---|---|---|
| | **Parallel** | **Non-Parallel** | |
| GEFCOM | 0:00:48 | 0:02:34 | 0.032 |
| Jena Weather | 0:06:11 | 0:21:19 | 0.007 |
| Household Energy | 0:37:58 | 1:30:24 | 0.013 |

TABLE III

EFFICIENCY AND ACCURACY OF THE 3CS ALGORITHM ON LARGE-SCALE DATASETS

# Conclusions and Future Work

- In this paper:

  - **We proposed the** *Concatenated Composite Covariance Search (3CS)* **algorithm**

  - **We evaluated that algorithm's capabilities by means of eight benchmark datasets and compared its performance to given state-of-the-art methods**

- As future work, we plan to incorporate global approximations into the procedure, develop domain-specific adaptations and include prior knowledge into the retrieval process

# Thanks for your attention!

**Fabian Berns**

Einsteinstraße 62
48149 Münster

Email: fabian.berns@uni-muenster.de
Website: www.uni-muenster.de/dma

**living.knowledge**