



Aggregating Dependent Gaussian Experts in Local Approximation

Hamed Jalali, Gjergji Kasneci

January 13, 2021



Gaussian Process (GP):

- Regression problem $y = f(x) + \epsilon$, where $x \in R^D$ and $\epsilon \sim \mathcal{N}(0, \sigma^2)$.
- Objective: learn f from a training set $\mathcal{D} = \{X,y\}_{i=1}^n$

-
$$f \sim GP\left(m(x), k(x, x')\right)$$

• Training: maximise the log-marginal likelihood

$$\log p(y|X) = -\frac{1}{2}y^{T}\mathcal{C}^{-1}y - \frac{1}{2}\log|\mathcal{C}| - \frac{n}{2}\log(2\pi)$$

where $\mathcal{C} = K + \sigma^2 I$.

• Cost: $\mathcal{O}(n^3)$ (because of the inversion and determinant of \mathcal{C})





Distributed GP reduces the cost of the standard GP by distributing the training process.

- \mathcal{D} is divided into M partitions $\mathcal{D}_1, \ldots, \mathcal{D}_M$, (called experts)
- Partitions are called experts
- The predictive distribution of the i'th expert \mathcal{M}_i and test input x^* is $p_i(y^*|\mathcal{D}_i, x^*) \sim \mathcal{N}(\mu_i^*, \Sigma_i^*)$,

$$\mu_i^* = k_{i*}^T (K_i + \sigma^2 I)^{-1} y_i,$$

$$\Sigma_i^* = k_{**} - k_{i*}^T (K_i + \sigma^2 I)^{-1} k_{i*}.$$





- The local Gaussian experts are conditionally independent.
 - CI assumption
- The predictive distribution of DGP is

$$p(y^*|\mathcal{D}, x^*) \propto \prod_{i=1}^M p_i^{\beta_i}(y^*|\mathcal{D}_i, x^*).$$

• The weights $\beta = \{\beta_1, \dots, \beta_M\}$ describe the importance of the experts.

The ensembles based on CI return sub-optimal solutions.



- Instead of the Gaussian experts $\mathcal{M} = \{\mathcal{M}_1, \dots, \mathcal{M}_M\}$, we define clusters of correlated experts, $\mathcal{C} = \{\mathcal{C}_1, \dots, \mathcal{C}_P\}$, $P \ll M$.
- Aggregating the experts at each cluster leads to a new layer of experts, K = {K₁,...,K_P}, which are conditionally independent given y*.









- let $\mu_{\mathcal{M}}^* = [\mu_1^*, \dots, \mu_M^*]^T$ be a $n_t \times M$ matrix that contains the local predictions of M experts at n_t test points.
- Assumption: The joint distribution of the experts predictions is multivariate normal.
- Gaussian graphical model:

$$p(\mu_{\mathcal{M}}^*|h,\Omega) \propto \exp\left\{-\frac{1}{2}(\mu_{\mathcal{M}}^*)^T \ \Omega \ \mu_{\mathcal{M}}^* + h^T \mu_{\mathcal{M}}^*\right\}.$$

- Ω is the precision matrix and encodes the conditional dependency.
- Ω is calculated using the GLasso method

$$\hat{\Omega} = \arg \max_{\Omega} \log |\Omega| - trace(S\Omega) - \lambda \|\Omega\|_{1}.$$





- The precision matrix is used to find the clusters of experts set \mathcal{C} .
- Each cluster C_i contains strongly dependent experts.
- We apply spectral clustering (SC) to find $C = \{C_1, \ldots, C_P\}$.
- SC makes use of the similarity matrix (here Ω).
- To aggregate the experts at each cluster, we use GRBCM method that leads to $\mathcal{K} = \{\mathcal{K}_1, \dots, \mathcal{K}_P\}.$
- New experts are conditionally independent.



Experts clustering









Experiments-Synthetic Example









	Pumadyn		Kin40k		Sacros		Song		
Model	SMSE	MSLL	SMSE	MSLL	SMSE	MSLL	SMSE	MSLL	_
DGEA (OURS) POE GPOE BCM RBCM GRBCM	0.0486 0.0505 0.0505 0.0499 0.0498 0.0511	-1.5133 4.8725 -1.4936 4.6688 12.1101 -1.488	0.0538 0.856 0.0856 0.0818 0.0772 0.0544	-1.3025 2.4153 -1.2286 1.6974 2.5256 -1.2785	0.0269 0.0311 0.0311 0.0308 0.0305 0.0305	-1.823 25.2807 -1.7756 24.868 61.5392 -1.4308	0.8084 0.8169 0.8169 10.4291 5.4373 0.8268	-0.122 69.9464 -0.123 44.1745 1.2089 0.2073	





We proposed a novel approach that leverages the dependencies between experts and improves the prediction quality. It

- uses an undirected graphical model to detect strong dependencies between experts
- defines clusters of interdependent experts
- provides consistent results when $n \to \infty$.

