

Scalable Deep Gaussian Markov Random Fields for General Graphs

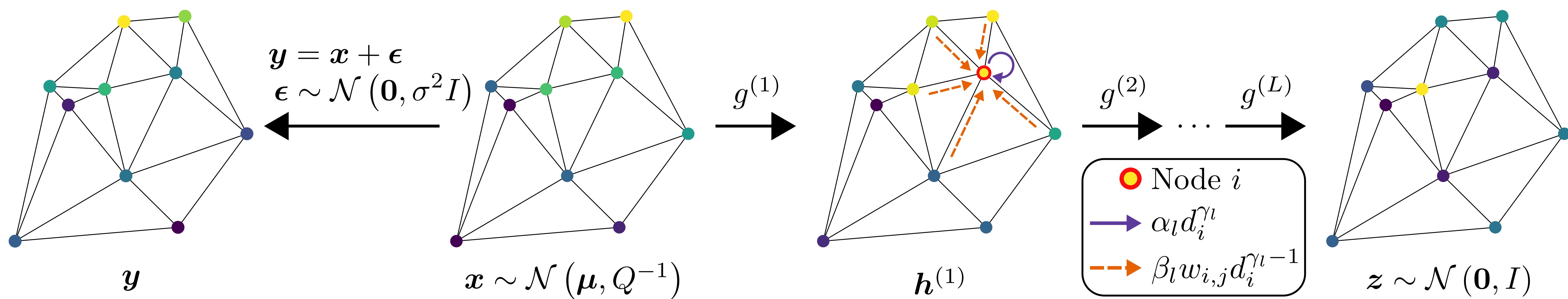
Joel Oskarsson¹

Per Sidén^{1,2}

Fredrik Lindsten¹

¹Division of Statistics and Machine Learning, Department of Computer and Information Science, Linköping University, Linköping, Sweden

²Arriver Software AB



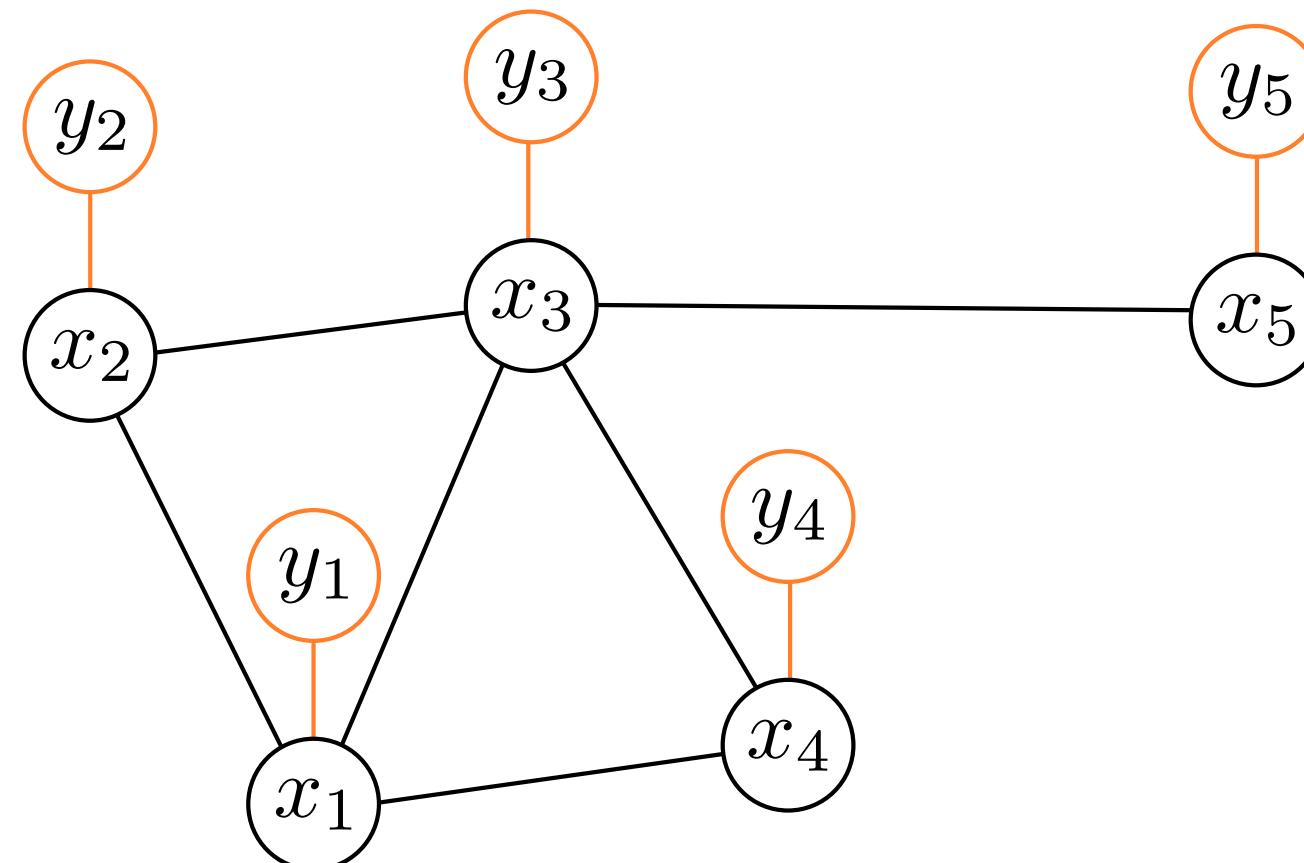
A Layered Gaussian Model on Graphs

We propose a flexible Gaussian Markov Random Field (GMRF) model for general graphs built on the multi-layer structure of Deep GMRFs. By designing a new type of layer we enable the model to scale to large graphs. The layer allows for efficient training using variational inference and existing software frameworks for Graph Neural Networks. For a Gaussian likelihood, close to exact Bayesian inference is available for the latent field. In experiments on synthetic and real world datasets our model compares favorably to other both Bayesian and deep learning methods.

Background: Gaussian Markov Random Fields

Gaussian Markov Random Fields (GMRFs)

- Graphical models where the nodes x are jointly Gaussian [2]
- Conjugate prior to Gaussian likelihood $p(y|x) = \mathcal{N}(y|x, \sigma^2 I)$
 - Posterior $p(x|y)$ analytically tractable, also a GMRF!



Deep GMRFs (DGMRFs) [3]

- GMRF x defined implicitly by affine map g
- $$z = g(x) = Gx + b, \quad z \sim \mathcal{N}(0, I) \quad (1)$$
- g defined as a combination of L simpler layers $g = g^{(L)} \circ g^{(L-1)} \circ \dots \circ g^{(1)}$
 - Originally restricted to lattice graphs (image-structured data)

Method: Graph DGMRF

- Consider graph \mathcal{G} with adjacency matrix A and degree matrix $D = \text{diag}([d_1, d_2, \dots, d_N]^\top)$
 - Define a layer construct
- $$\mathbf{h}^{(l)} = g^{(l)}(\mathbf{h}^{(l-1)}) = G^{(l)}\mathbf{h}^{(l-1)} + b_l \mathbf{1} \quad (2)$$
- $$G^{(l)} = \alpha_l D^{\gamma_l} + \beta_l D^{\gamma_l - 1} A \quad (3)$$
- b_l, α_l, β_l and γ_l are trainable parameters
 - Closely corresponds to Graph Neural Network (GNN) formulation. Existing frameworks can be utilized for automatic differentiation and GPU acceleration.
 - Variational training, maximizing the Evidence Lower Bound (ELBO)
 - Training requires efficient evaluation of $\log |\det(G^{(l)})|$. We develop two solutions:
 - Exact method based on pre-computing the eigenvalues of $D^{-1}A$
 - Scalable approximate method based on reformulating log-determinant as a power series [1]

Experiments

- Experiments on prediction for a subset of unobserved nodes. We train on a fraction of the nodes and evaluate on remaining ones.

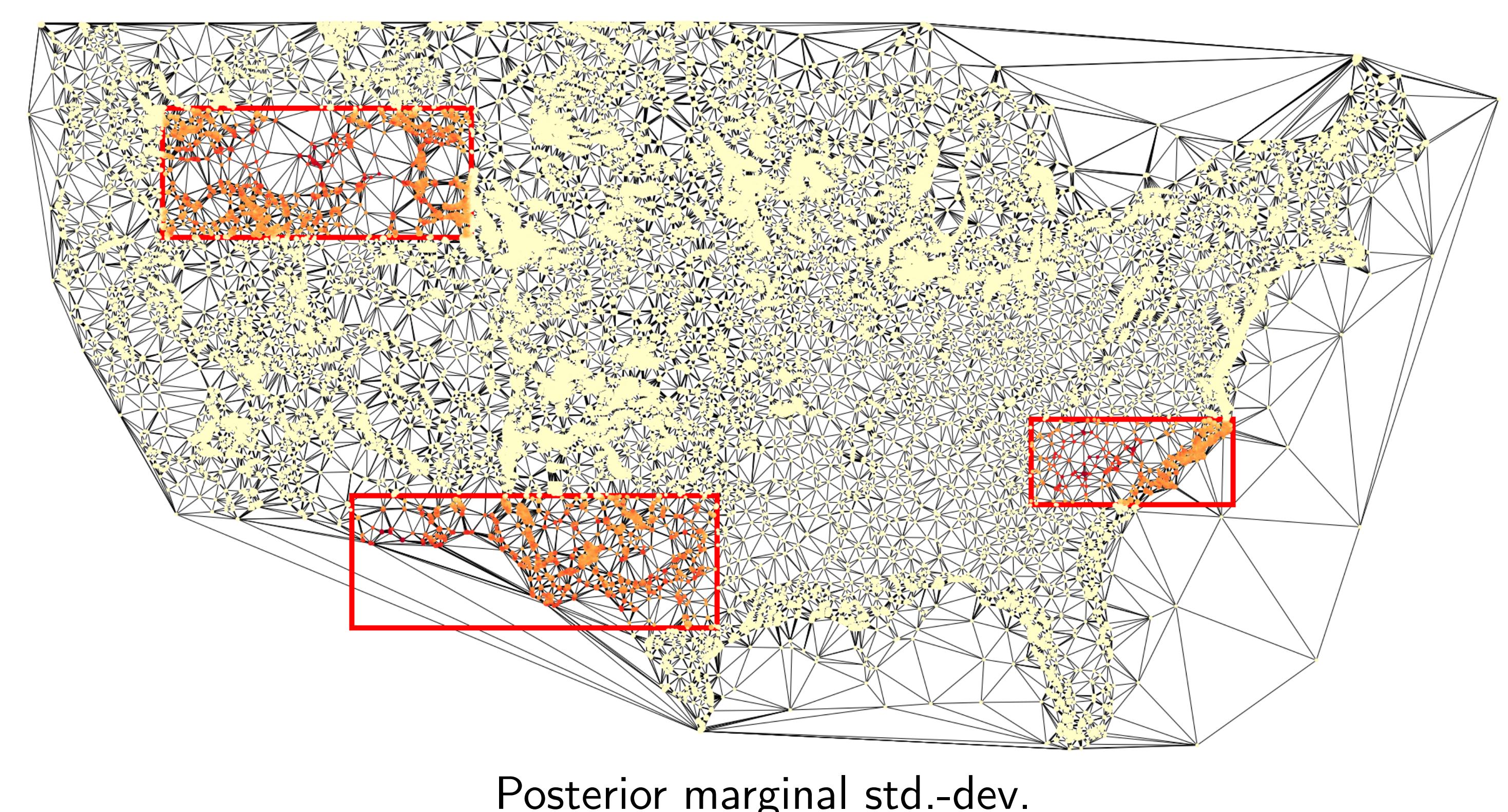
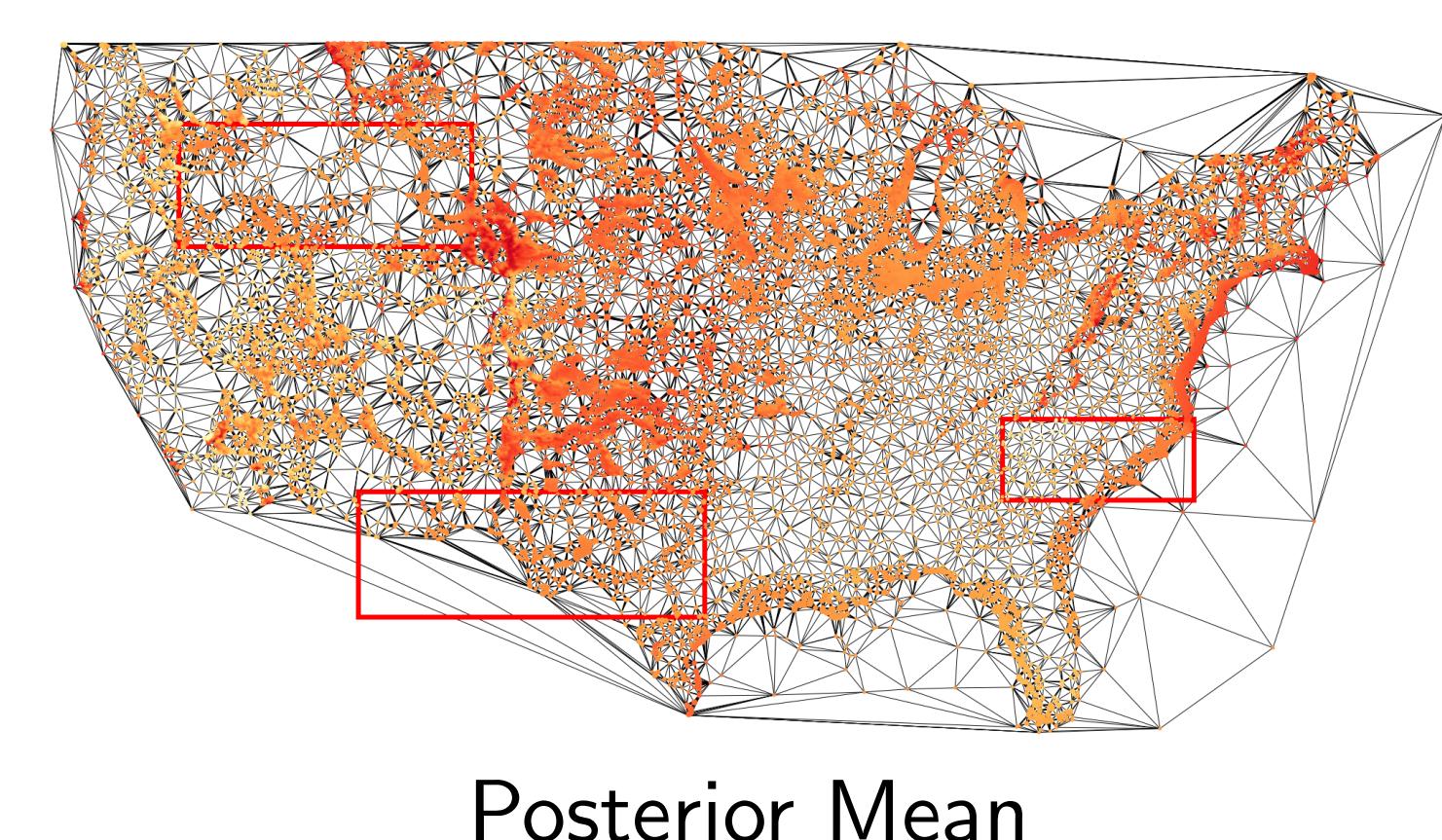
Wikipedia Graph

- 11 631 Wikipedia pages as nodes
- Edges are links between pages
- Target y is log of average monthly traffic
- CRPS, probabilistic metric taking uncertainty into account

	RMSE	CRPS
Baselines		
GRAPH GP	2.169	1.251
DGP (GNN)	1.308	0.786
I-GMRF	1.526	0.939
DGMRF		
1 LAYER	1.311	0.704
3 LAYER	1.228	0.652
5 LAYER	1.169	0.614

Wind Speed Dataset

- Spatial graph with 126 652 nodes
- Target y is average wind speed
- Nodes inside red rectangles are unobserved



More Information



Code, Link to Paper:

github.com/joeloskarsson/graph-dgmrf

Correspondence to:

Joel Oskarsson, joel.oskarsson@liu.se

References

- J. Behrmann, W. Grathwohl, et al. Invertible residual networks. In *Proceedings of the 36th International Conference on Machine Learning*, 2019.
- H. Rue and L. Held. *Gaussian Markov random fields: theory and applications*. Number 104 in Monographs on statistics and applied probability. Chapman & Hall/CRC, 2005.
- P. Sidén and F. Lindsten. Deep gaussian markov random fields. In *Proceedings of the 37th International Conference on Machine Learning*, 2020.

Scalable Deep Gaussian Markov Random Fields for General Graphs

Joel Oskarsson¹

Per Sidén^{1,2}

Fredrik Lindsten¹

¹Division of Statistics and Machine Learning, Department of Computer and Information Science, Linköping University, Linköping, Sweden
²Arriver Software AB

Variational Training

- Directly maximizing the log marginal likelihood $\log p(\mathbf{y}_m|\theta)$ is not computationally feasible for large graphs
- Maximizing the ELBO allows for scalable and efficient learning of the model parameters
- Let $\mathbf{m} \in \{0, 1\}^N$ be an observation mask with 1 at observed nodes

$$M = \sum_{i=1}^N m_i \quad \mathbf{y}_m = \mathbf{y} \odot \mathbf{m} \quad I_m = \text{diag}(\mathbf{m})$$

$$\begin{aligned} \text{ELBO}(\theta, \phi) = & -\frac{1}{2} \mathbb{E}_{q(\mathbf{x}|\phi)} \left[g(\mathbf{x})^\top g(\mathbf{x}) + \frac{1}{\sigma^2} (\mathbf{y}_m - \mathbf{x})^\top I_m (\mathbf{y}_m - \mathbf{x}) \right] \\ & + \sum_{l=1}^L \log |\det(G^{(l)})| + H[q(\mathbf{x}|\phi)] - M \log \sigma + \text{const.} \end{aligned} \quad (4)$$

Variational Distribution

- Gaussian variational distribution $q(\mathbf{x}|\phi) = \mathcal{N}(\mathbf{x}|\boldsymbol{\nu}, SS^\top)$
- Defined in opposite direction of DGMRF

$$\mathbf{x} = S\mathbf{r} + \boldsymbol{\nu}, \quad \mathbf{r} \sim \mathcal{N}(\mathbf{0}, I) \quad (5)$$

- [3] use mean-field approximation (diagonal S)
- We can re-use our layer construct through $S = \text{diag}(\boldsymbol{\xi}) \tilde{G} \text{diag}(\boldsymbol{\tau})$, where \tilde{G} is one or more layers as defined in Eq. 3.
- Introduces off-diagonal elements in the covariance matrix of q

Computing the Log-Determinant

Eigenvalue Method

$$\log |\det(G^{(l)})| = \sum_{i=1}^N \gamma_l \log(d_i) + \log |\alpha_l + \beta_l \lambda'_i| \quad (6)$$

- $\{\lambda'_i\}_{i=1}^N$ are the eigenvalues of $D^{-1}A$
 - Depends only on graph structure \Rightarrow Can be pre-computed!
- Exact, but does not scale to massive graphs

Power Series Method

- With $\tilde{A} = D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$

$$\log |\det(G^{(l)})| = N \log(\alpha_l) + \sum_{i=1}^N \gamma_l \log(d_i) + \log \left| \det \left(I + \frac{\beta_l}{\alpha_l} \tilde{A} \right) \right| \quad (7)$$

- Formulate last log-determinant as power series [1]

$$\log \left| \det \left(I + \frac{\beta_l}{\alpha_l} \tilde{A} \right) \right| = \sum_{k=1}^{\infty} -\frac{1}{k} \left(-\frac{\beta_l}{\alpha_l} \right)^k \text{Tr}(\tilde{A}^k). \quad (8)$$

- Truncate at some large $k = K$ during training

– $\text{Tr}(\tilde{A}^k)$ can be pre-computed for $k \in \{1, \dots, K\}$

- Approximate, but highly scalable

Acknowledgements

This research is financially supported by the Swedish Research Council via the project *Handling Uncertainty in Machine Learning Systems* (contract number: 2020-04122), the Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by the Knut and Alice Wallenberg Foundation, and the Excellence Center at Linköping–Lund in Information Technology (ELLIIT).

Posterior Inference

- GMRF prior is conjugate to Gaussian likelihood $p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{x}, \sigma^2 I)$

Prior

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, Q^{-1})$$

$$\boldsymbol{\mu} = -G^{-1}\mathbf{b}$$

$$Q = G^\top G$$

$$\mathbf{x} | \mathbf{y}_m \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{Q}^{-1})$$

$$\tilde{\boldsymbol{\mu}} = \tilde{Q}^{-1} \left(Q\boldsymbol{\mu} + \frac{1}{\sigma^2} \mathbf{y}_m \right)$$

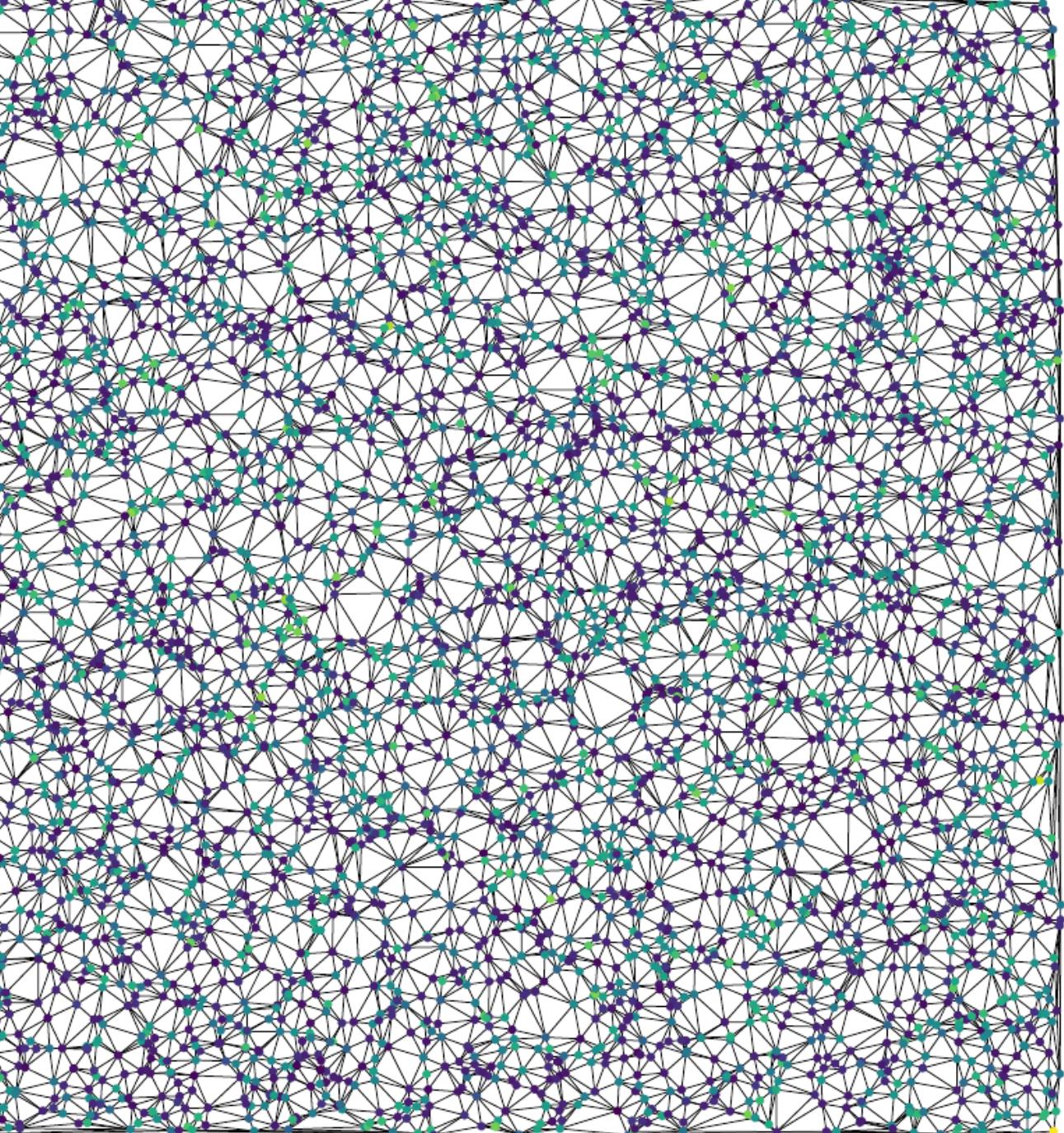
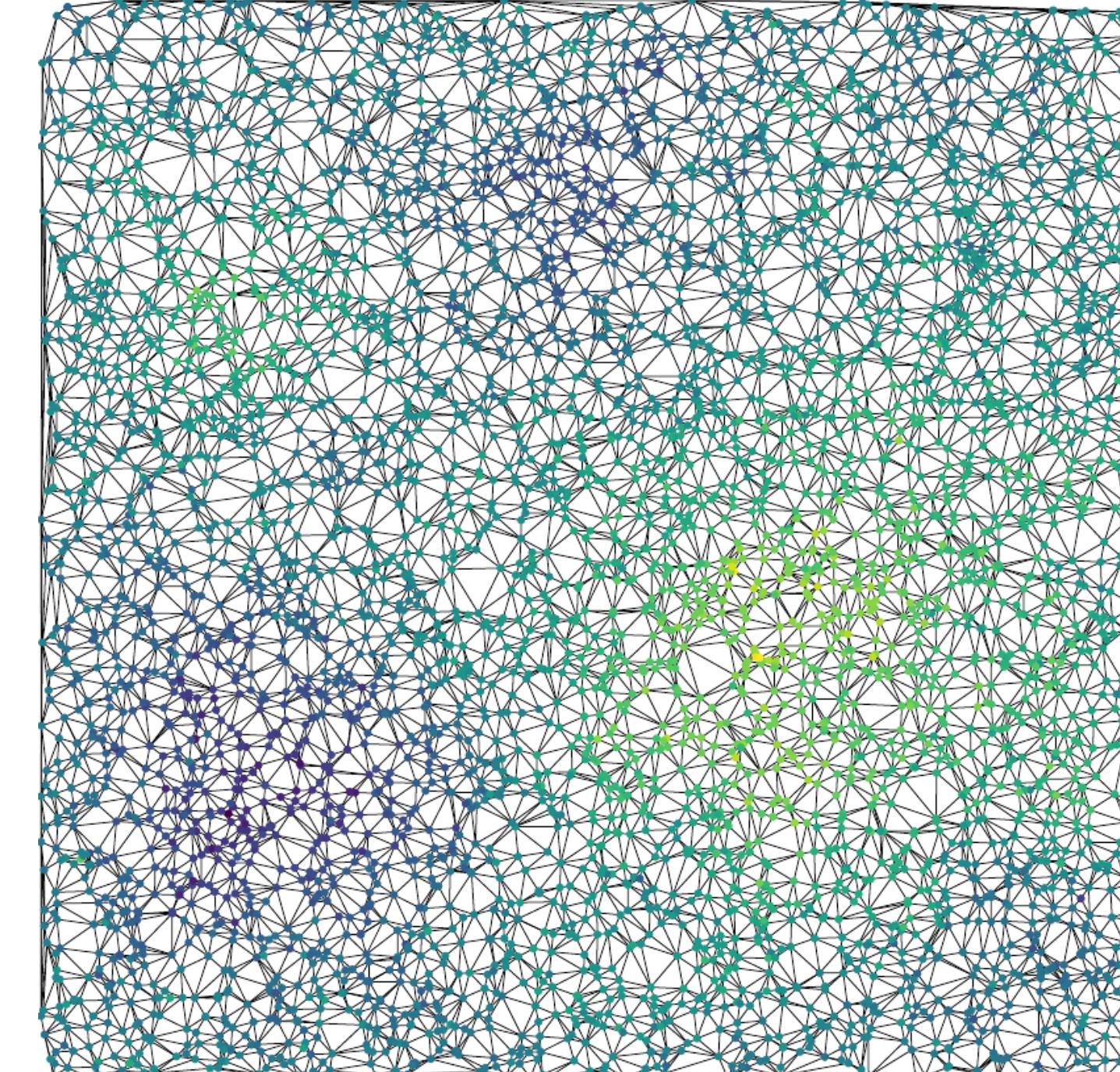
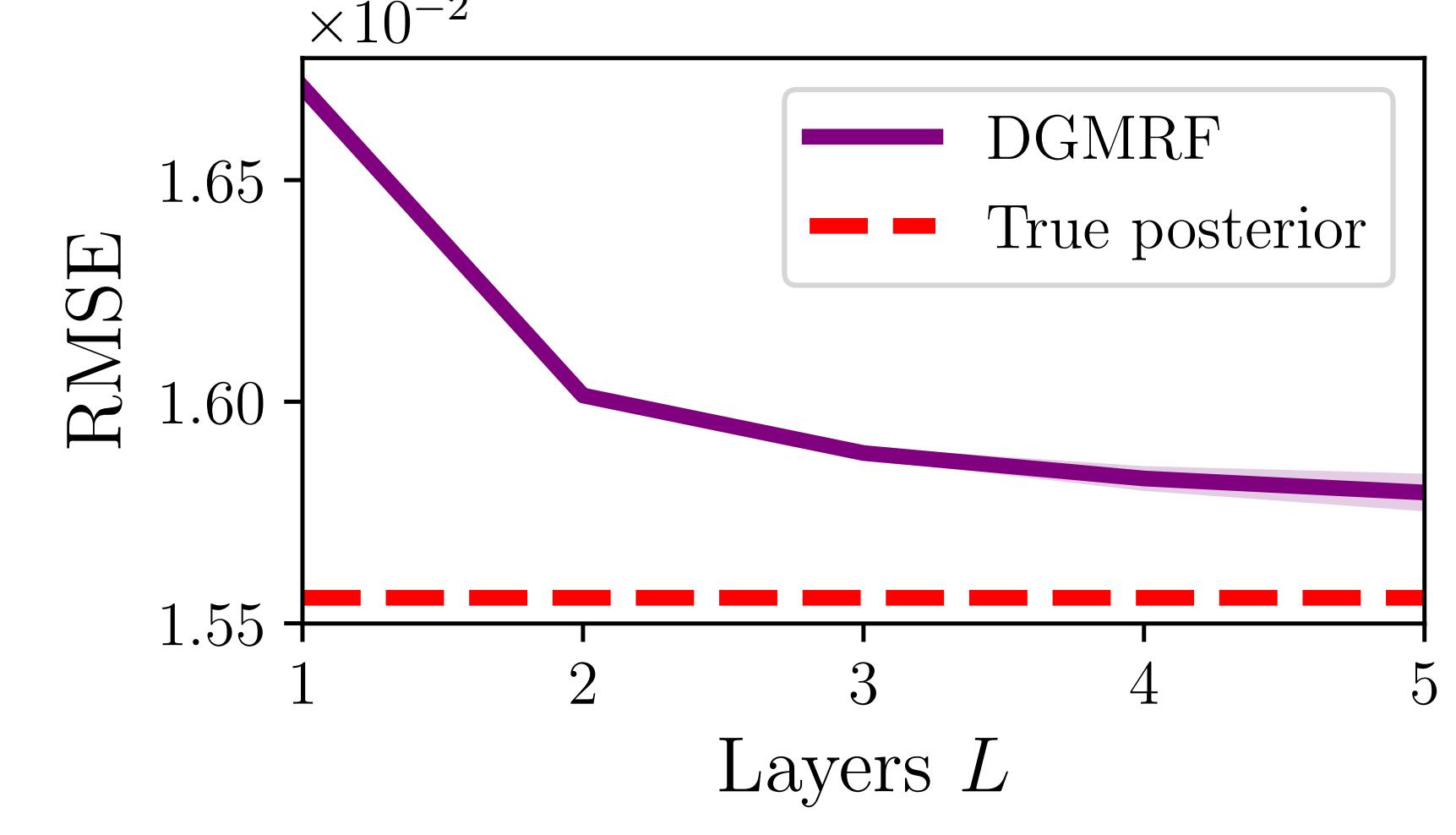
$$\tilde{Q} = Q + \frac{1}{\sigma^2} I_m$$

- Conjugate Gradient method used to avoid computing inverse in \tilde{Q}

Additional Experiments

Synthetic GMRF Dataset

- Sample from known GMRF created as random mix of multiple DGMRFs
- 1-5 layer DGMRFs trained and evaluated



California Housing Dataset

- Data on 20 640 housing blocks with socio-economic features
- Target \mathbf{y} is log of median house value in block
- Spatial graph with edges weighted by inverse distances
- We consider the dataset both with and without node features
 - In **No features** the methods use only the graph/spatial coordinates
 - In DGMRFs node features are handled using an auxiliary Bayesian linear model

Model	No features		Features	
	RMSE	CRPS	RMSE	CRPS
BAYES LR	13.319	7.521	8.872	4.834
MLP	11.086	7.915	7.094	4.525
GCN	8.760	5.683	6.837	4.273
GAT	9.166	6.049	6.788	4.348
SVGP	10.172	5.689	7.287	3.930
GRAPH GP	11.202	6.350	-	-
IGMRF	6.989	3.841	-	-
DGMRF, $L = 1$	6.909	3.665	5.894	3.078
DGMRF, $L = 2$	6.853	3.651	5.810	3.041
DGMRF, $L = 3$	6.853	3.656	5.804	3.039