

Scalable Bayesian inference for Latent Gaussian Models



King Abdullah University of
Science and Technology



جامعة الملك عبد الله
للعلوم والتقنية

Håvard Rue
King Abdullah University of Science and Technology
Saudi Arabia

June 2026



- Importance of GMRFs
- Latent Gaussian Models
- Bayesian Inference
- Key component: (a new) Sparse matrix solver



- Importance of GMRFs
- Latent Gaussian Models
- Bayesian Inference
- Key component: (a new) Sparse matrix solver



- Importance of GMRFs
- Latent Gaussian Models
- Bayesian Inference
- Key component: (a new) Sparse matrix solver



- Importance of GMRFs
- Latent Gaussian Models
- Bayesian Inference
- Key component: (a new) Sparse matrix solver

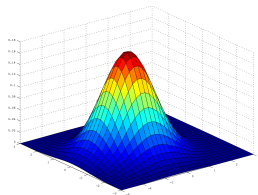


Multivariate Normal distribution

- Random vector (x_1, \dots, x_n)
- (Scaled and centred) density

$$\pi(\mathbf{x}) \propto |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}\right)$$

- $\boldsymbol{\Sigma}$ is the (SPD) covariance matrix
- $\Sigma_{ij} = E(x_i x_j)$
- $\Sigma_{ij} = 0 \Leftrightarrow x_i \perp x_j (i \neq j)$



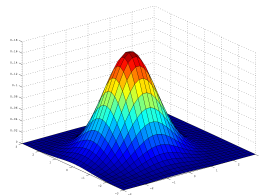


Multivariate Normal distribution

- Random vector (x_1, \dots, x_n)
- (Scaled and centred) density

$$\pi(\mathbf{x}) \propto |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}\right)$$

- $\boldsymbol{\Sigma}$ is the (SPD) covariance matrix
- $\Sigma_{ij} = E(x_i x_j)$
- $\Sigma_{ij} = 0 \Leftrightarrow x_i \perp x_j (i \neq j)$



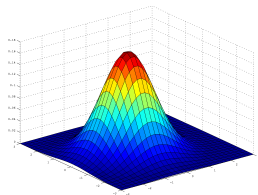


Multivariate Normal distribution

- Random vector (x_1, \dots, x_n)
- (Scaled and centred) density

$$\pi(\mathbf{x}) \propto |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}\right)$$

- $\boldsymbol{\Sigma}$ is the (SPD) covariance matrix
- $\Sigma_{ij} = E(x_i x_j)$
- $\Sigma_{ij} = 0 \Leftrightarrow x_i \perp x_j (i \neq j)$



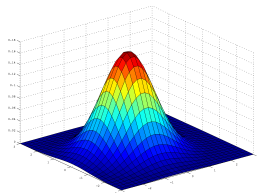


Multivariate Normal distribution

- Random vector (x_1, \dots, x_n)
- (Scaled and centred) density

$$\pi(\mathbf{x}) \propto |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}\right)$$

- $\boldsymbol{\Sigma}$ is the (SPD) covariance matrix
- $\Sigma_{ij} = E(x_i x_j)$
- $\Sigma_{ij} = 0 \Leftrightarrow x_i \perp x_j (i \neq j)$



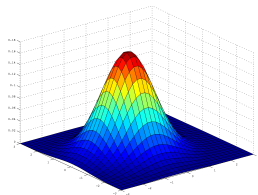


Multivariate Normal distribution

- Random vector (x_1, \dots, x_n)
- (Scaled and centred) density

$$\pi(\mathbf{x}) \propto |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}\right)$$

- $\boldsymbol{\Sigma}$ is the (SPD) covariance matrix
- $\Sigma_{ij} = E(x_i x_j)$
- $\Sigma_{ij} = 0 \Leftrightarrow x_i \perp x_j (i \neq j)$





Conditional independence and GMRFs

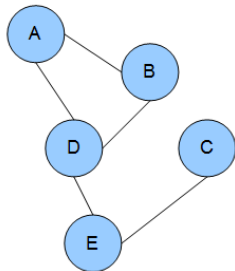
- Conditional independence of x_i and x_j , equals independence in

$$\pi(x_i, x_j \mid \mathbf{x}_{-ij})$$

- The result is, with $\mathbf{Q} = \Sigma^{-1}$, that

$$Q_{ij} = 0 \Leftrightarrow x_i \perp x_j \mid \mathbf{x}_{-ij}$$

- This is potentially very very useful, as
 - If \mathbf{Q} is fully connected then Σ is dense (Cayley-Hamilton thm)
 - Explain dependence by a sparse \mathbf{Q}
 - Sparse matrices: faster computations





Conditional independence and GMRFs

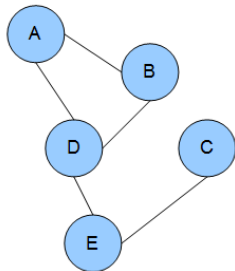
- Conditional independence of x_i and x_j , equals independence in

$$\pi(x_i, x_j \mid \mathbf{x}_{-ij})$$

- The result is, with $\mathbf{Q} = \Sigma^{-1}$, that

$$Q_{ij} = 0 \Leftrightarrow x_i \perp x_j \mid \mathbf{x}_{-ij}$$

- This is potentially very very useful, as
 - If \mathbf{Q} is fully connected then Σ is dense (Cayley-Hamilton thm)
 - Explain dependence by a sparse \mathbf{Q}
 - Sparse matrices: faster computations





Conditional independence and GMRFs

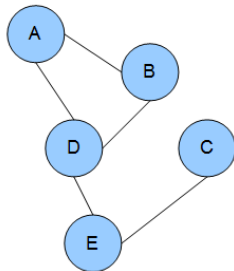
- Conditional independence of x_i and x_j , equals independence in

$$\pi(x_i, x_j \mid \mathbf{x}_{-ij})$$

- The result is, with $\mathbf{Q} = \Sigma^{-1}$, that

$$Q_{ij} = 0 \Leftrightarrow x_i \perp x_j \mid \mathbf{x}_{-ij}$$

- This is potentially very very useful, as
 - If \mathbf{Q} is fully connected then Σ is dense (Cayley-Hamilton thm)
 - Explain dependence by a sparse \mathbf{Q}
 - Sparse matrices: faster computations





Conditional independence and GMRFs

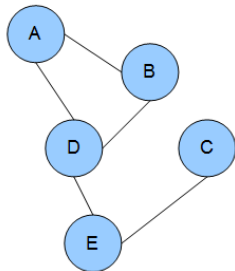
- Conditional independence of x_i and x_j , equals independence in

$$\pi(x_i, x_j \mid \mathbf{x}_{-ij})$$

- The result is, with $\mathbf{Q} = \boldsymbol{\Sigma}^{-1}$, that

$$Q_{ij} = 0 \Leftrightarrow x_i \perp x_j \mid \mathbf{x}_{-ij}$$

- This is potentially very very useful, as
 - If \mathbf{Q} is fully connected then $\boldsymbol{\Sigma}$ is dense (Cayley–Hamilton thm)
 - Explain dependence by a sparse \mathbf{Q}
 - Sparse matrices: faster computations





Conditional independence and GMRFs

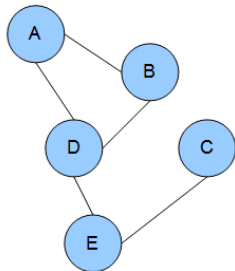
- Conditional independence of x_i and x_j , equals independence in

$$\pi(x_i, x_j \mid \mathbf{x}_{-ij})$$

- The result is, with $\mathbf{Q} = \boldsymbol{\Sigma}^{-1}$, that

$$Q_{ij} = 0 \Leftrightarrow x_i \perp x_j \mid \mathbf{x}_{-ij}$$

- This is potentially very very useful, as
 - If \mathbf{Q} is fully connected then $\boldsymbol{\Sigma}$ is dense (Cayley–Hamilton thm)
 - Explain dependence by a sparse \mathbf{Q}
 - Sparse matrices: faster computations





Conditional independence and GMRFs

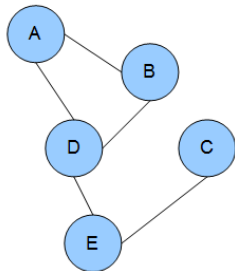
- Conditional independence of x_i and x_j , equals independence in

$$\pi(x_i, x_j \mid \mathbf{x}_{-ij})$$

- The result is, with $\mathbf{Q} = \boldsymbol{\Sigma}^{-1}$, that

$$Q_{ij} = 0 \Leftrightarrow x_i \perp x_j \mid \mathbf{x}_{-ij}$$

- This is potentially very very useful, as
 - If \mathbf{Q} is fully connected then $\boldsymbol{\Sigma}$ is dense (Cayley–Hamilton thm)
 - Explain dependence by a sparse \mathbf{Q}
 - Sparse matrices: faster computations





Example: Markov processes in time

- Auto-regressive process of order 1

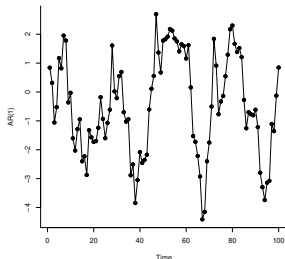
$$x_t = \phi x_{t-1} + \epsilon_t, \quad 0 \leq \phi < 1, \quad t = 1, \dots, n$$

with Gaussian noise ϵ_t

- Linear graph
- Q is tridiagonal (factorisation cost $\mathcal{O}(n)$)
- Σ is $\phi^{-|i-j|}$ (exp-correlation function)
- Cont.time: Ornstein-Uhlenbeck process

$$dx_t = -\theta x_t dt + dW_t,$$

Wiener process W_t





Example: Markov processes in time

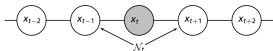
- Auto-regressive process of order 1

$$x_t = \phi x_{t-1} + \epsilon_t, \quad 0 \leq \phi < 1, \quad t = 1, \dots, n$$

with Gaussian noise ϵ_t

- Linear graph

- Q is tridiagonal (factorisation cost $\mathcal{O}(n)$)
- Σ is $\phi^{-|i-j|}$ (exp-correlation function)
- Cont.time: Ornstein-Uhlenbeck process



$$dx_t = -\theta x_t dt + dW_t,$$

Wiener process W_t



Example: Markov processes in time

- Auto-regressive process of order 1

$$x_t = \phi x_{t-1} + \epsilon_t, \quad 0 \leq \phi < 1, \quad t = 1, \dots, n$$

with Gaussian noise ϵ_t

- Linear graph
- \mathbf{Q} is tridiagonal (factorisation cost $\mathcal{O}(n)$)
- Σ is $\phi^{-|i-j|}$ (exp-correlation function)
- Cont.time: Ornstein-Uhlenbeck process

$$\begin{pmatrix} 1 & \phi & \phi^2 & \phi^3 & \phi^4 & \phi^5 & \phi^6 \\ \phi & 1 & \phi & \phi^2 & \phi^3 & \phi^4 & \phi^5 \\ \phi^2 & \phi & 1 & \phi & \phi^2 & \phi^3 & \phi^4 \\ \phi^3 & \phi^2 & \phi & 1 & \phi & \phi^2 & \phi^3 \\ \phi^4 & \phi^3 & \phi^2 & \phi & 1 & \phi & \phi^2 \\ \phi^5 & \phi^4 & \phi^3 & \phi^2 & \phi & 1 & \phi \\ \phi^6 & \phi^5 & \phi^4 & \phi^3 & \phi^2 & \phi & 1 \end{pmatrix}$$

$$dx_t = -\theta x_t dt + dW_t,$$

Wiener process W_t



Example: Markov processes in time

- Auto-regressive process of order 1

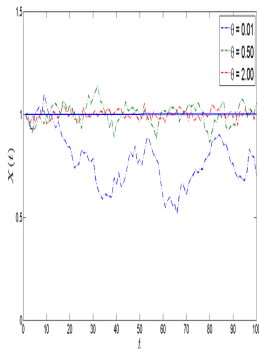
$$x_t = \phi x_{t-1} + \epsilon_t, \quad 0 \leq \phi < 1, \quad t = 1, \dots, n$$

with Gaussian noise ϵ_t

- Linear graph
- \mathbf{Q} is tridiagonal (factorisation cost $\mathcal{O}(n)$)
- Σ is $\phi^{-|i-j|}$ (exp-correlation function)
- Cont.time: Ornstein-Uhlenbeck process

$$dx_t = -\theta x_t dt + dW_t,$$

Wiener process W_t





How does correlation work?

- Let

$$x_1 = z_1 + \mu$$

$$x_2 = z_2 + \mu$$

- z_1, z_2, μ are independent zero mean Gaussian's
- $\text{Var}(\mu) = \sigma^2, \text{Var}(z_i) = 1$
- Then

$$\text{Corr}(x_1, x_2) = \frac{\sigma^2}{\sigma^2 + 1}$$

- This is a constructive/intuitive way to introduce correlation (which is a 'consequence' of the construction)



How does correlation work?

- Let

$$x_1 = z_1 + \mu$$

$$x_2 = z_2 + \mu$$

- z_1, z_2, μ are independent zero mean Gaussian's
- $\text{Var}(\mu) = \sigma^2, \text{Var}(z_i) = 1$
- Then

$$\text{Corr}(x_1, x_2) = \frac{\sigma^2}{\sigma^2 + 1}$$

- This is an constructive/intuitive way to introduce correlation (which is a 'consequence' of the construction)



How does correlation work?

- Let

$$x_1 = z_1 + \mu$$

$$x_2 = z_2 + \mu$$

- z_1, z_2, μ are independent zero mean Gaussian's
- $\text{Var}(\mu) = \sigma^2, \text{Var}(z_i) = 1$
- Then

$$\text{Corr}(x_1, x_2) = \frac{\sigma^2}{\sigma^2 + 1}$$

- This is an constructive/intuitive way to introduce correlation (which is a 'consequence' of the construction)



How does correlation work?

- Let

$$x_1 = z_1 + \mu$$

$$x_2 = z_2 + \mu$$

- z_1, z_2, μ are independent zero mean Gaussian's
- $\text{Var}(\mu) = \sigma^2, \text{Var}(z_i) = 1$
- Then

$$\text{Corr}(x_1, x_2) = \frac{\sigma^2}{\sigma^2 + 1}$$

- This is an constructive/intuitive way to introduce correlation (which is a 'consequence' of the construction)



How does correlation work?

- Let

$$x_1 = z_1 + \mu$$

$$x_2 = z_2 + \mu$$

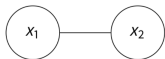
- z_1, z_2, μ are independent zero mean Gaussian's
- $\text{Var}(\mu) = \sigma^2, \text{Var}(z_i) = 1$
- Then

$$\text{Corr}(x_1, x_2) = \frac{\sigma^2}{\sigma^2 + 1}$$

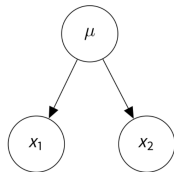
- This is an constructive/intuitive way to introduce correlation (which is a 'consequence' of the construction)



How this influence sparsity?



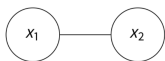
$$e_{x_1, x_2} = \begin{pmatrix} \times & \times \\ \times & \times \end{pmatrix}$$



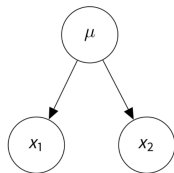
$$e_{x_1, x_2, \mu} = \begin{pmatrix} \times & & \times \\ & \times & \times \\ \times & \times & \times \end{pmatrix}$$



How this influence sparsity?



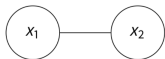
$$\mathbf{e}_{x_1, x_2} = \begin{pmatrix} \times & \times \\ \times & \times \end{pmatrix}$$



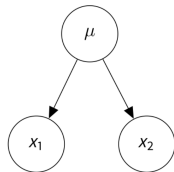
$$\mathbf{e}_{x_1, x_2, \mu} = \begin{pmatrix} \times & & \times \\ & \times & \times \\ \times & \times & \times \end{pmatrix}$$



How this influence sparsity?



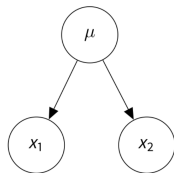
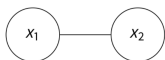
$$\mathbf{e}_{x_1, x_2} = \begin{pmatrix} \times & \times \\ \times & \times \end{pmatrix}$$



$$\mathbf{e}_{x_1, x_2, \mu} = \begin{pmatrix} \times & & \times \\ & \times & \times \\ \times & \times & \times \end{pmatrix}$$



How this influence sparsity?



$$\mathbf{e}_{x_1, x_2} = \begin{pmatrix} \times & \times \\ \times & \times \end{pmatrix}$$

$$\mathbf{e}_{x_1, x_2, \mu} = \begin{pmatrix} \times & & \times \\ & \times & \times \\ \times & \times & \times \end{pmatrix}$$



What are good “operations” for GMRFs?

Conditioning $\mathbf{x}_1 | \mathbf{x}_2$

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12}^T & \Sigma_{22} \end{pmatrix}$$

$$\mathbf{Q} = \begin{pmatrix} \mathbf{Q}_{11} & \mathbf{Q}_{12} \\ \mathbf{Q}_{12}^T & \mathbf{Q}_{22} \end{pmatrix}$$

$$\text{Var}(\mathbf{x}_1 | \mathbf{x}_2) = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{12}^T$$

$$\text{Prec}(\mathbf{x}_1 | \mathbf{x}_2) = \mathbf{Q}_{11}$$



What are good “operations” for GMRFs?

Updating/Learning:

$$\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_x^{-1})$$

then we observe

$$\mathbf{y}|\mathbf{x} \sim \mathcal{N}(\mathbf{x}, \mathbf{Q}_y^{-1})$$

$$\text{Var}(\mathbf{x}|\mathbf{y}) = (\boldsymbol{\Sigma}_x^{-1} + \boldsymbol{\Sigma}_y^{-1})^{-1}$$

$$\text{Prec}(\mathbf{x}|\mathbf{y}) = \mathbf{Q}_x + \mathbf{Q}_y$$



What are good “operations” for GMRFs?

Updating/Learning:

$$\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_x^{-1})$$

then we observe

$$\mathbf{y}|\mathbf{x} \sim \mathcal{N}(\mathbf{x}, \mathbf{Q}_y^{-1})$$

$$\text{Var}(\mathbf{x}|\mathbf{y}) = (\boldsymbol{\Sigma}_x^{-1} + \boldsymbol{\Sigma}_y^{-1})^{-1}$$

$$\text{Prec}(\mathbf{x}|\mathbf{y}) = \mathbf{Q}_x + \mathbf{Q}_y$$



What are good “operations” for GMRFs?

Updating/Learning:

$$\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_x^{-1})$$

then we observe

$$\mathbf{y}|\mathbf{x} \sim \mathcal{N}(\mathbf{x}, \mathbf{Q}_y^{-1})$$

$$\text{Var}(\mathbf{x}|\mathbf{y}) = (\boldsymbol{\Sigma}_x^{-1} + \boldsymbol{\Sigma}_y^{-1})^{-1}$$

$$\text{Prec}(\mathbf{x}|\mathbf{y}) = \mathbf{Q}_x + \mathbf{Q}_y$$



What are good “operations” for GMRFs?

Distribution of sums:

$$\boldsymbol{\eta} = \mathbf{x} + \mathbf{y} + \mathbf{z}$$

where $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_x^{-1})$, $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_y^{-1})$ and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_z^{-1})$.

$$\mathbf{Q}_\eta = (\mathbf{Q}_x^{-1} + \mathbf{Q}_y^{-1} + \mathbf{Q}_z^{-1})^{-1}$$

\mathbf{Q}_η is (in general) dense and it is costly to update \mathbf{Q}_η

The covariance formulation is easy

$$\boldsymbol{\Sigma}_\eta = \boldsymbol{\Sigma}_x + \boldsymbol{\Sigma}_y + \boldsymbol{\Sigma}_z$$

Marginalisation destroy sparsity...



What are good “operations” for GMRFs?

Distribution of sums:

$$\boldsymbol{\eta} = \mathbf{x} + \mathbf{y} + \mathbf{z}$$

where $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_x^{-1})$, $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_y^{-1})$ and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_z^{-1})$.

$$\mathbf{Q}_\eta = (\mathbf{Q}_x^{-1} + \mathbf{Q}_y^{-1} + \mathbf{Q}_z^{-1})^{-1}$$

\mathbf{Q}_η is (in general) dense and it is costly to update \mathbf{Q}_η

The covariance formulation is easy

$$\boldsymbol{\Sigma}_\eta = \boldsymbol{\Sigma}_x + \boldsymbol{\Sigma}_y + \boldsymbol{\Sigma}_z$$

Marginalisation destroy sparsity...



What are good “operations” for GMRFs?

Distribution of sums:

$$\boldsymbol{\eta} = \mathbf{x} + \mathbf{y} + \mathbf{z}$$

where $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_x^{-1})$, $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_y^{-1})$ and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_z^{-1})$.

$$\mathbf{Q}_\eta = (\mathbf{Q}_x^{-1} + \mathbf{Q}_y^{-1} + \mathbf{Q}_z^{-1})^{-1}$$

\mathbf{Q}_η is (in general) dense and it is costly to update \mathbf{Q}_η

The covariance formulation is easy

$$\boldsymbol{\Sigma}_\eta = \boldsymbol{\Sigma}_x + \boldsymbol{\Sigma}_y + \boldsymbol{\Sigma}_z$$

Marginalisation destroy sparsity...



What are good “operations” for GMRFs?

Distribution of sums:

$$\boldsymbol{\eta} = \boldsymbol{x} + \boldsymbol{y} + \boldsymbol{z}$$

where $\boldsymbol{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_x^{-1})$, $\boldsymbol{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_y^{-1})$ and $\boldsymbol{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_z^{-1})$.

$$\mathbf{Q}_\eta = (\mathbf{Q}_x^{-1} + \mathbf{Q}_y^{-1} + \mathbf{Q}_z^{-1})^{-1}$$

\mathbf{Q}_η is (in general) dense and it is costly to update \mathbf{Q}_η

The covariance formulation is easy

$$\boldsymbol{\Sigma}_\eta = \boldsymbol{\Sigma}_x + \boldsymbol{\Sigma}_y + \boldsymbol{\Sigma}_z$$

Marginalisation destroy sparsity...



What are good “operations” for GMRFs?

Distribution of sums:

$$\boldsymbol{\eta} = \mathbf{x} + \mathbf{y} + \mathbf{z}$$

where $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_x^{-1})$, $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_y^{-1})$ and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_z^{-1})$.

- Using “ ϵ ” with high precision τ

$$\boldsymbol{\eta} = \mathbf{x} + \mathbf{y} + \mathbf{z} + \boldsymbol{\epsilon}$$

- Then

$$\text{Var} \begin{pmatrix} \boldsymbol{\eta} \\ \mathbf{x} \\ \mathbf{y} \\ \mathbf{z} \end{pmatrix} = \begin{pmatrix} \tau I & -\tau I & -\tau I & -\tau I \\ & \mathbf{Q}_x + \tau I & \tau I & \tau I \\ & & \mathbf{Q}_y + \tau I & \tau I \\ & & & \mathbf{Q}_z + \tau I \end{pmatrix}$$

- Good practical solution unless choice of τ becomes an issue



What are good “operations” for GMRFs?

Distribution of sums:

$$\boldsymbol{\eta} = \mathbf{x} + \mathbf{y} + \mathbf{z}$$

where $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_x^{-1})$, $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_y^{-1})$ and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_z^{-1})$.

- Using “ ϵ ” with high precision τ

$$\boldsymbol{\eta} = \mathbf{x} + \mathbf{y} + \mathbf{z} + \boldsymbol{\epsilon}$$

- Then

$$\text{Var} \begin{pmatrix} \boldsymbol{\eta} \\ \mathbf{x} \\ \mathbf{y} \\ \mathbf{z} \end{pmatrix} = \begin{pmatrix} \tau \mathbf{I} & -\tau \mathbf{I} & -\tau \mathbf{I} & -\tau \mathbf{I} \\ & \mathbf{Q}_x + \tau \mathbf{I} & \tau \mathbf{I} & \tau \mathbf{I} \\ & & \mathbf{Q}_y + \tau \mathbf{I} & \tau \mathbf{I} \\ & & & \mathbf{Q}_z + \tau \mathbf{I} \end{pmatrix}$$

- Good practical solution unless choice of τ becomes an issue



What are good “operations” for GMRFs?

Distribution of sums:

$$\boldsymbol{\eta} = \mathbf{x} + \mathbf{y} + \mathbf{z}$$

where $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_x^{-1})$, $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_y^{-1})$ and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_z^{-1})$.

- Using “ ϵ ” with high precision τ

$$\boldsymbol{\eta} = \mathbf{x} + \mathbf{y} + \mathbf{z} + \boldsymbol{\epsilon}$$

- Then

$$\text{Var} \begin{pmatrix} \boldsymbol{\eta} \\ \mathbf{x} \\ \mathbf{y} \\ \mathbf{z} \end{pmatrix} = \begin{pmatrix} \tau \mathbf{I} & -\tau \mathbf{I} & -\tau \mathbf{I} & -\tau \mathbf{I} \\ & \mathbf{Q}_x + \tau \mathbf{I} & \tau \mathbf{I} & \tau \mathbf{I} \\ & & \mathbf{Q}_y + \tau \mathbf{I} & \tau \mathbf{I} \\ & & & \mathbf{Q}_z + \tau \mathbf{I} \end{pmatrix}$$

- Good practical solution unless choice of τ becomes an issue



What are good “operations” for GMRFs?

Distribution of sums:

$$\boldsymbol{\eta} = \mathbf{x} + \mathbf{y} + \mathbf{z}$$

where $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_x^{-1})$, $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_y^{-1})$ and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_z^{-1})$.

- Using cumulative sums

$$\begin{pmatrix} \tilde{\mathbf{x}} \\ \tilde{\mathbf{y}} \\ \tilde{\mathbf{z}} \end{pmatrix} = \begin{pmatrix} \mathbf{x} + \mathbf{y} + \mathbf{z} \\ \mathbf{y} + \mathbf{z} \\ \mathbf{z} \end{pmatrix}$$

- Then

$$\text{Var} \begin{pmatrix} \tilde{\mathbf{x}} \\ \tilde{\mathbf{y}} \\ \tilde{\mathbf{z}} \end{pmatrix} = \begin{pmatrix} \mathbf{Q}_x & & & \\ & -\mathbf{Q}_x & & \\ & & \mathbf{Q}_y + \mathbf{Q}_x & -\mathbf{Q}_y \\ & & & \mathbf{Q}_z + \mathbf{Q}_y \end{pmatrix}$$

- More efficient, but no access to \mathbf{x} and \mathbf{y} directly.



What are good “operations” for GMRFs?

Distribution of sums:

$$\boldsymbol{\eta} = \mathbf{x} + \mathbf{y} + \mathbf{z}$$

where $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_x^{-1})$, $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_y^{-1})$ and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_z^{-1})$.

- Using cumulative sums

$$\begin{pmatrix} \tilde{\mathbf{x}} \\ \tilde{\mathbf{y}} \\ \tilde{\mathbf{z}} \end{pmatrix} = \begin{pmatrix} \mathbf{x} + \mathbf{y} + \mathbf{z} \\ \mathbf{y} + \mathbf{z} \\ \mathbf{z} \end{pmatrix}$$

- Then

$$\text{Var} \begin{pmatrix} \tilde{\mathbf{x}} \\ \tilde{\mathbf{y}} \\ \tilde{\mathbf{z}} \end{pmatrix} = \begin{pmatrix} \mathbf{Q}_x & & & \\ & -\mathbf{Q}_x & & \\ & & \dots & \\ & & \dots & -\mathbf{Q}_y \\ & & & & \dots & \\ & & & & & \dots & \mathbf{Q}_z + \mathbf{Q}_y \end{pmatrix}$$

- More efficient, but no access to \mathbf{x} and \mathbf{y} directly.



What are good “operations” for GMRFs?

Distribution of sums:

$$\boldsymbol{\eta} = \mathbf{x} + \mathbf{y} + \mathbf{z}$$

where $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_x^{-1})$, $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_y^{-1})$ and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_z^{-1})$.

- Using cumulative sums

$$\begin{pmatrix} \tilde{\mathbf{x}} \\ \tilde{\mathbf{y}} \\ \tilde{\mathbf{z}} \end{pmatrix} = \begin{pmatrix} \mathbf{x} + \mathbf{y} + \mathbf{z} \\ \mathbf{y} + \mathbf{z} \\ \mathbf{z} \end{pmatrix}$$

- Then

$$\text{Var} \begin{pmatrix} \tilde{\mathbf{x}} \\ \tilde{\mathbf{y}} \\ \tilde{\mathbf{z}} \end{pmatrix} = \begin{pmatrix} \mathbf{Q}_x & & & \\ & -\mathbf{Q}_x & & \\ & & \dots & \\ & & \dots & -\mathbf{Q}_y \\ & & & & \dots & \\ & & & & & \dots & \mathbf{Q}_z + \mathbf{Q}_y \end{pmatrix}$$

- More efficient, but no access to \mathbf{x} and \mathbf{y} directly.



Additive regression models/GLM++

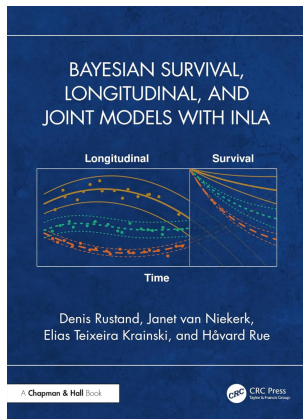
- Linear predictor

$$\eta_i = \mu + \sum_j \beta_j z_{ij} + \mathbf{f}_{1i} + \mathbf{f}_{2i} + \dots$$

with observations

$$y_i | \dots \sim \pi(y_i | \eta_i, \dots)$$

- \mathbf{f}_k is a Gaussian process, think times-series, spatial-model, spline, some kind of dependent ‘random effect’, measurement error model, etc...





Additive regression models/GLM++

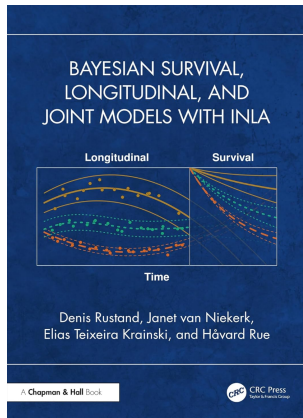
- Linear predictor

$$\eta_i = \mu + \sum_j \beta_j z_{ij} + \mathbf{f}_{1i} + \mathbf{f}_{2i} + \dots$$

with observations

$$y_i | \dots \sim \pi(y_i | \eta_i, \dots)$$

- \mathbf{f}_k is a Gaussian process, think times-series, spatial-model, spline, some kind of dependent ‘random effect’, measurement error model, etc...





Key observation

- If each \mathbf{f}_k is a GMRF, then

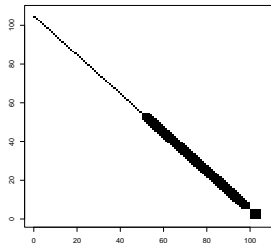
$$\mathbf{x} = (\beta_1, \beta_2, \dots, \mathbf{f}_1, \mathbf{f}_2, \dots)$$

is also a GMRF with precision matrix \mathbf{Q}

- Let $\boldsymbol{\eta} = \mathbf{A}\mathbf{x}$. Assume Gaussian data with unit variance, then $\mathbf{x}|\mathbf{y}$ is a GMRF with precision matrix

$$\mathbf{Q} + \mathbf{A}^T \mathbf{A}$$

- Non-Gaussian data: “sequence of Gaussian-data”





Key observation

- If each \mathbf{f}_k is a GMRF, then

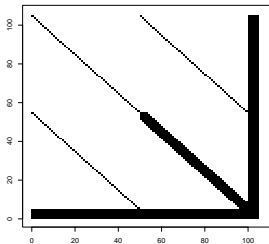
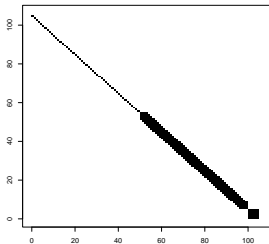
$$\mathbf{x} = (\beta_1, \beta_2, \dots, \mathbf{f}_1, \mathbf{f}_2, \dots)$$

is also a GMRF with precision matrix \mathbf{Q}

- Let $\boldsymbol{\eta} = \mathbf{A}\mathbf{x}$. Assume Gaussian data with unit variance, then $\mathbf{x}|\mathbf{y}$ is a GMRF with precision matrix

$$\mathbf{Q} + \mathbf{A}^T \mathbf{A}$$

- Non-Gaussian data: “sequence of Gaussian-data”





Key observation

- If each \mathbf{f}_k is a GMRF, then

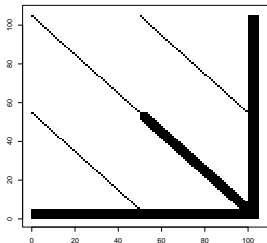
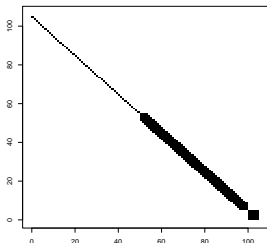
$$\mathbf{x} = (\beta_1, \beta_2, \dots, \mathbf{f}_1, \mathbf{f}_2, \dots)$$

is also a GMRF with precision matrix \mathbf{Q}

- Let $\boldsymbol{\eta} = \mathbf{A}\mathbf{x}$. Assume Gaussian data with unit variance, then $\mathbf{x}|\mathbf{y}$ is a GMRF with precision matrix

$$\mathbf{Q} + \mathbf{A}^T \mathbf{A}$$

- Non-Gaussian data: “sequence of Gaussian-data”





Keep it in the family!

GMRFs is closed

- sums
- conditioning

iff done right: no marginalisation!





Keep it in the family!

GMRFs is closed

- sums
- conditioning

iff done right: no marginalisation!





Keep it in the family!

GMRFs is closed

- sums
- conditioning

iff done right: no marginalisation!





Keep it in the family!

GMRFs is closed

- sums
- conditioning

iff done right: no marginalisation!





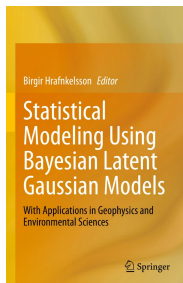
Bayesian inference: Latent Gaussian models (LGMs)

The most important class of models...

- Conditionally independent data $y_i | \eta_i, \theta$
- Latent

$$\eta_i = \sum_j \beta_j c_{ij} + \sum_k f_{ki}$$

- Hyperparameters θ (latent and/or likelihood)
- Each component is a GMRF (or a “small” dense)
- Covariance is a property, not a defining quantity





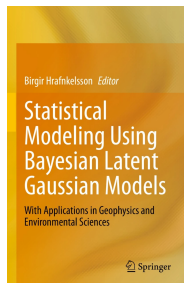
Bayesian inference: Latent Gaussian models (LGMs)

The most important class of models...

- Conditionally independent data $y_i | \eta_i, \theta$
- Latent

$$\eta_i = \sum_j \beta_j c_{ij} + \sum_k f_{ki}$$

- Hyperparameters θ (latent and/or likelihood)
- Each component is a GMRF (or a “small” dense)
- Covariance is a property, not a defining quantity





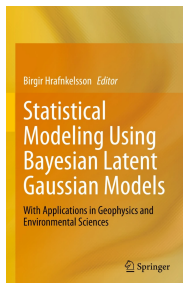
Bayesian inference: Latent Gaussian models (LGMs)

The most important class of models...

- Conditionally independent data $y_i | \eta_i, \theta$
- Latent

$$\eta_i = \sum_j \beta_j c_{ij} + \sum_k f_{ki}$$

- Hyperparameters θ (latent and/or likelihood)
- Each component is a GMRF (or a “small” dense)
- Covariance is a property, not a defining quantity





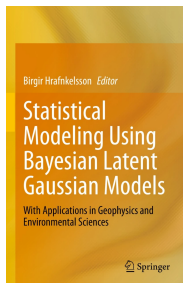
Bayesian inference: Latent Gaussian models (LGMs)

The most important class of models...

- Conditionally independent data $y_i | \eta_i, \theta$
- Latent

$$\eta_i = \sum_j \beta_j c_{ij} + \sum_k f_{ki}$$

- Hyperparameters θ (latent and/or likelihood)
- Each component is a GMRF (or a “small” dense)
- Covariance is a property, not a defining quantity





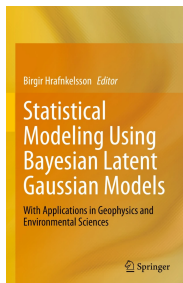
Bayesian inference: Latent Gaussian models (LGMs)

The most important class of models...

- Conditionally independent data $y_i | \eta_i, \theta$
- Latent

$$\eta_i = \sum_j \beta_j c_{ij} + \sum_k f_{ki}$$

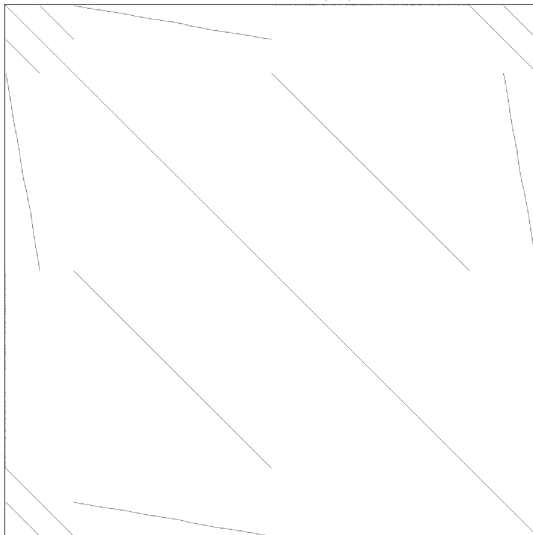
- Hyperparameters θ (latent and/or likelihood)
- Each component is a GMRF (or a “small” dense)
- Covariance is a property, not a defining quantity





Examples

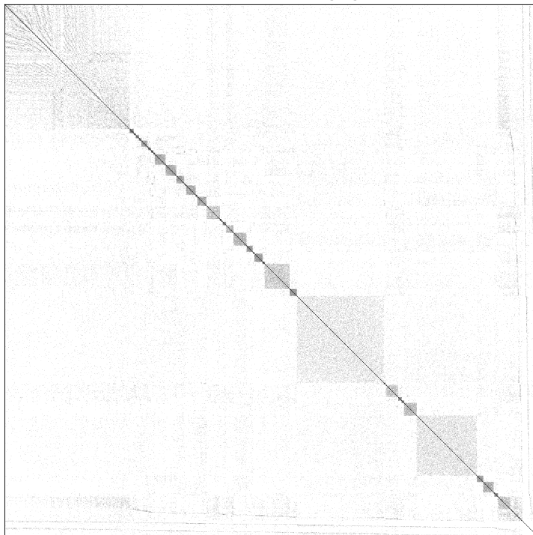
inla_graph 8rtKSK
n = 7,860 nnz = 58,670 avg deg = 7.4





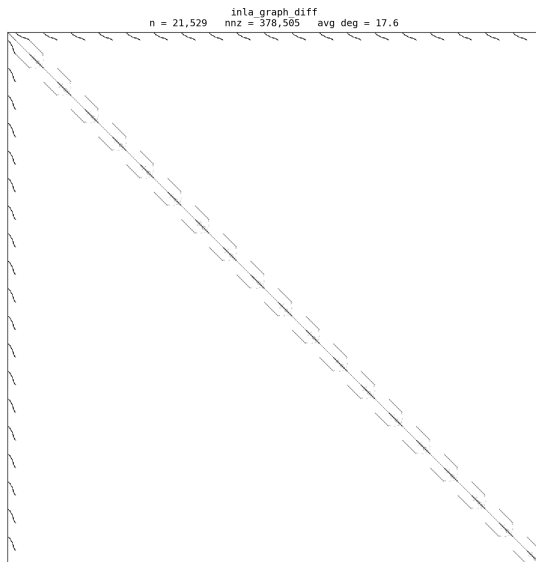
Examples

`inla_graph_ayaLrw`
n = 69,224 nnz = 267,944 avg deg = 3.9



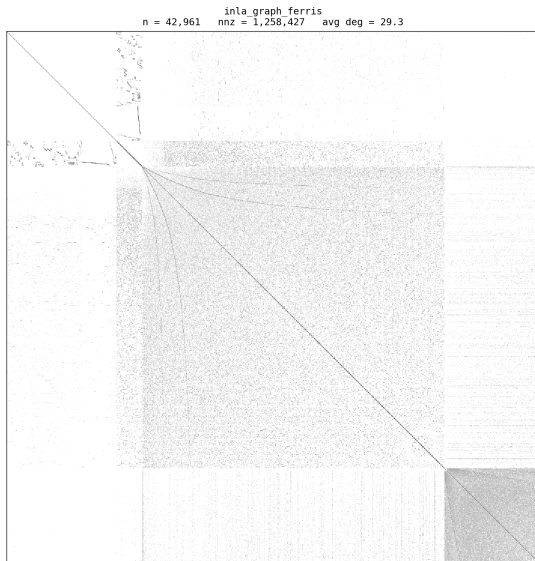


Examples





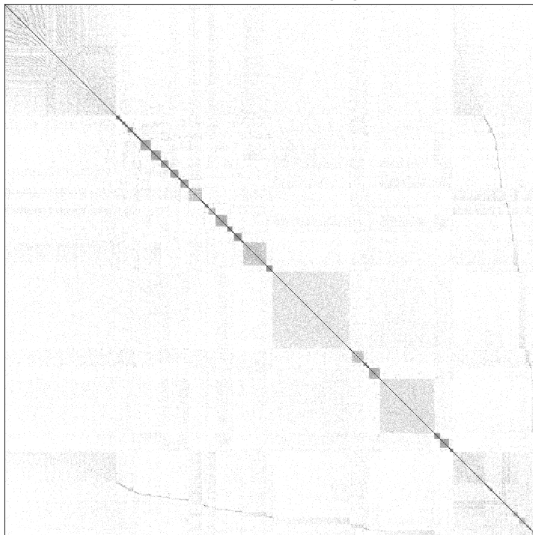
Examples





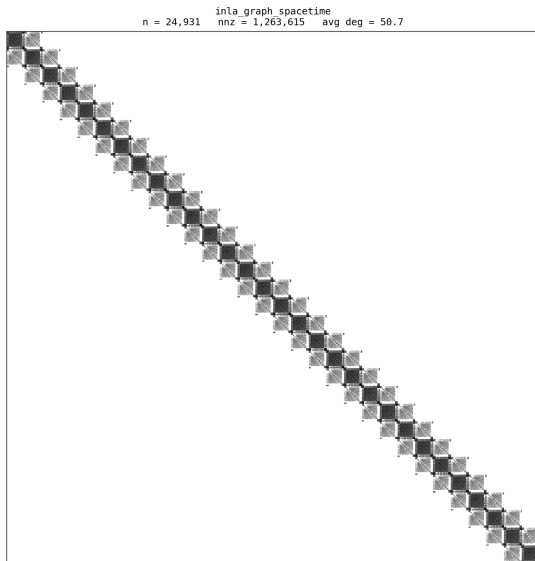
Examples

inla_graph_net814381
n = 81,438 nnz = 297,626 avg deg = 3.7



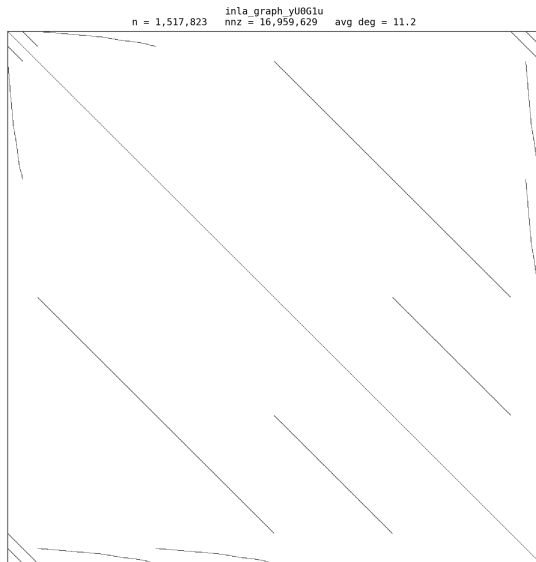


Examples





Examples





Bayesian Inference: simple ideas

- Gaussian approximations
- Laplace approximations
- Variational Bayes corrections
- Numerical integration
- Available in R-INLA package (www.r-inla.org)



Bayesian Inference: simple ideas

- Gaussian approximations
- Laplace approximations
- Variational Bayes corrections
- Numerical integration
- Available in R-INLA package (www.r-inla.org)



Bayesian Inference: simple ideas

- Gaussian approximations
- Laplace approximations
- Variational Bayes corrections
- Numerical integration
- Available in R-INLA package (www.r-inla.org)



Bayesian Inference: simple ideas

- Gaussian approximations
- Laplace approximations
- Variational Bayes corrections
- Numerical integration
- Available in R-INLA package (www.r-inla.org)



Bayesian Inference: simple ideas

- Gaussian approximations
- Laplace approximations
- Variational Bayes corrections
- Numerical integration
- Available in R-INLA package (www.r-inla.org)



Computational issues

Numerical methods for sparse matrices

- Factorise $Q = LL^T$
- Partial inverse Q^{-1} for all $i \sim j$ or $i = j$.
- Determinant $\log |Q|$
- Solve $Qx = b$ and $QX = B$

```

for (i = 0; i < dimensionSize; i++) {
  for (j = 0; j <= i; j++) {
    float sum = 0;
    for (k = 0; k < j; k++)
      sum += L[i][k] * L[j][k];

    if (i == j)
      L[i][j] = sqrt(A[i][i] - sum);
    else
      L[i][j] = (1.0 / L[j][j]) * (A[i][j] - sum);
  }
}

```

Most important ingredient to scalable inference, is to handle sparse matrices well!



Computational issues

Numerical methods for sparse matrices

- Factorise $\mathbf{Q} = \mathbf{L}\mathbf{L}^T$
- Partial inverse \mathbf{Q}^{-1} for all $i \sim j$ or $i = j$.
- Determinant $\log |\mathbf{Q}|$
- Solve $\mathbf{Q}\mathbf{x} = \mathbf{b}$ and $\mathbf{Q}\mathbf{X} = \mathbf{B}$

```

for (i = 0; i < dimensionSize; i++) {
  for (j = 0; j <= i; j++) {
    float sum = 0;
    for (k = 0; k < j; k++)
      sum += L[i][k] * L[j][k];

    if (i == j)
      L[i][j] = sqrt(A[i][i] - sum);
    else
      L[i][j] = (1.0 / L[j][j]) * (A[i][j] - sum);
  }
}

```

Most important ingredient to scalable inference, is to handle sparse matrices well!



Computational issues

Numerical methods for sparse matrices

- Factorise $\mathbf{Q} = \mathbf{L}\mathbf{L}^T$
- Partial inverse \mathbf{Q}^{-1} for all $i \sim j$ or $i = j$.
- Determinant $\log |\mathbf{Q}|$
- Solve $\mathbf{Q}\mathbf{x} = \mathbf{b}$ and $\mathbf{Q}\mathbf{X} = \mathbf{B}$

```

for (i = 0; i < dimensionSize; i++) {
  for (j = 0; j <= i; j++) {
    float sum = 0;
    for (k = 0; k < j; k++)
      sum += L[i][k] * L[j][k];

    if (i == j)
      L[i][j] = sqrt(A[i][i] - sum);
    else
      L[i][j] = (1.0 / L[j][j]) * (A[i][j] - sum);
  }
}

```

Most important ingredient to scalable inference, is to handle sparse matrices well!



Computational issues

Numerical methods for sparse matrices

- Factorise $\mathbf{Q} = \mathbf{L}\mathbf{L}^T$
- Partial inverse \mathbf{Q}^{-1} for all $i \sim j$ or $i = j$.
- Determinant $\log |\mathbf{Q}|$
- Solve $\mathbf{Q}\mathbf{x} = \mathbf{b}$ and $\mathbf{Q}\mathbf{X} = \mathbf{B}$

```

for (i = 0; i < dimensionSize; i++) {
  for (j = 0; j <= i; j++) {
    float sum = 0;
    for (k = 0; k < j; k++)
      sum += L[i][k] * L[j][k];

    if (i == j)
      L[i][j] = sqrt(A[i][i] - sum);
    else
      L[i][j] = (1.0 / L[j][j]) * (A[i][j] - sum);
  }
}

```

Most important ingredient to scalable inference, is to handle sparse matrices well!



Computational issues

Numerical methods for sparse matrices

- Factorise $\mathbf{Q} = \mathbf{L}\mathbf{L}^T$
- Partial inverse \mathbf{Q}^{-1} for all $i \sim j$ or $i = j$.
- Determinant $\log |\mathbf{Q}|$
- Solve $\mathbf{Q}\mathbf{x} = \mathbf{b}$ and $\mathbf{Q}\mathbf{X} = \mathbf{B}$

```

for (i = 0; i < dimensionSize; i++) {
  for (j = 0; j <= i; j++) {
    float sum = 0;
    for (k = 0; k < j; k++)
      sum += L[i][k] * L[j][k];

    if (i == j)
      L[i][j] = sqrt(A[i][i] - sum);
    else
      L[i][j] = (1.0 / L[j][j]) * (A[i][j] - sum);
  }
}

```

Most important ingredient to scalable inference, is to handle sparse matrices well!



Computational issues

Numerical methods for sparse matrices

- Factorise $\mathbf{Q} = \mathbf{L}\mathbf{L}^T$
- Partial inverse \mathbf{Q}^{-1} for all $i \sim j$ or $i = j$.
- Determinant $\log |\mathbf{Q}|$
- Solve $\mathbf{Q}\mathbf{x} = \mathbf{b}$ and $\mathbf{Q}\mathbf{X} = \mathbf{B}$

```

for (i = 0; i < dimensionSize; i++) {
  for (j = 0; j <= i; j++) {
    float sum = 0;
    for (k = 0; k < j; k++)
      sum += L[i][k] * L[j][k];

    if (i == j)
      L[i][j] = sqrt(A[i][i] - sum);
    else
      L[i][j] = (1.0 / L[j][j]) * (A[i][j] - sum);
  }
}

```

Most important ingredient to scalable inference, is to handle sparse matrices well!



Not that much around...

- HPC community like “dense” more than “sparse”
- **R-INLA** uses TAUCS (outdated > 20 years ago)
- CHOLMOD (what is used in R, MATLAB, does not scale well in parallel and do not handle inverse)
- PARDISO (good, not open-source and is now commercial, so...)
- etc...
- With the right person around, we can “just write our own”, right?



Esmail Abdul Fattah



Not that much around...

- HPC community like “dense” more than “sparse”
- **R-INLA** uses TAUCS (outdated > 20 years ago)
- CHOLMOD (what is used in R, MATLAB, does not scale well in parallel and do not handle inverse)
- PARDISO (good, not open-source and is now commercial, so...)
- etc...
- With the right person around, we can “just write our own”, right?



Esmail Abdul Fattah



Not that much around...

- HPC community like “dense” more than “sparse”
- **R-INLA** uses TAUCS (outdated > 20 years ago)
- CHOLMOD (what is used in R, MATLAB, does not scale well in parallel and do not handle inverse)
- PARDISO (good, not open-source and is now commercial, so...)
- etc...
- With the right person around, we can “just write our own”, right?



Esmail Abdul Fattah



Not that much around...

- HPC community like “dense” more than “sparse”
- **R-INLA** uses TAUCS (outdated > 20 years ago)
- CHOLMOD (what is used in R, MATLAB, does not scale well in parallel and do not handle inverse)
- PARDISO (good, not open-source and is now commercial, so...)
- etc...
- With the right person around, we can “just write our own”, right?



Esmail Abdul Fattah



Not that much around...

- HPC community like “dense” more than “sparse”
- **R-INLA** uses TAUCS (outdated > 20 years ago)
- CHOLMOD (what is used in R, MATLAB, does not scale well in parallel and do not handle inverse)
- PARDISO (good, not open-source and is now commercial, so...)
- etc...
- With the right person around, we can “just write our own”, right?



Esmail Abdul Fattah



Not that much around...

- HPC community like “dense” more than “sparse”
- **R-INLA** uses TAUCS (outdated > 20 years ago)
- CHOLMOD (what is used in R, MATLAB, does not scale well in parallel and do not handle inverse)
- PARDISO (good, not open-source and is now commercial, so...)
- etc...
- With the right person around, we can “just write our own”, right?



Esmail Abdul Fattah



Cholesky factorisation

```
for (i = 0; i < dimensionSize; i++) {
  for (j = 0; j <= i; j++) {
    float sum = 0;
    for (k = 0; k < j; k++)
      sum += L[i][k] * L[j][k];

    if (i == j)
      L[i][j] = sqrt(A[i][i] - sum);
    else
      L[i][j] = (1.0 / L[j][j] * (A[i][j] - sum));
  }
}
```




3 years later...

sparse • smart • structured • scalable

A High-Performance Framework for Sparse Matrices

sTiles targets the full **Cholesky** → **solve** → **selected inverse** pipeline on a single shared-memory node, using tile-based parallelism and GPU acceleration. Today's release handles symmetric positive-definite matrices across the entire spectrum, from very sparse to fully dense, with one unified solver. Distributed-memory support is on the roadmap.

[View Examples](#) [API Documentation](#)



Sparse Semi-Sparse Semi-Dense Dense





3 years later...

- We're getting there
- Announced in Glasgow'25 in May
- ...the we found some new problematic cases
- that took another year to fix
- Written for modern CPU and will have GPU acceleration (for hard cases)
- Results are *very* good, with excellent parallel scaling
- Handle hard cases particularly well
- Open source when done
- Stay tuned!





3 years later...

- We're getting there
- Announced in Glasgow'25 in May
- ...the we found some new problematic cases
- that took another year to fix
- Written for modern CPU and will have GPU acceleration (for hard cases)
- Results are *very* good, with excellent parallel scaling
- Handle hard cases particularly well
- Open source when done
- Stay tuned!





3 years later...

- We're getting there
- Announced in Glasgow'25 in May
- ...the we found some new problematic cases
- that took another year to fix
- Written for modern CPU and will have GPU acceleration (for hard cases)
- Results are *very* good, with excellent parallel scaling
- Handle hard cases particularly well
- Open source when done
- Stay tuned!





3 years later...

- We're getting there
- Announced in Glasgow'25 in May
- ...the we found some new problematic cases
- that took another year to fix
- Written for modern CPU and will have GPU acceleration (for hard cases)
- Results are *very good*, with excellent parallel scaling
- Handle hard cases particularly well
- Open source when done
- Stay tuned!





3 years later...

- We're getting there
- Announced in Glasgow'25 in May
- ...the we found some new problematic cases
- that took another year to fix
- Written for modern CPU and will have GPU acceleration (for hard cases)
- Results are *very good*, with excellent parallel scaling
- Handle hard cases particularly well
- Open source when done
- Stay tuned!





3 years later...

- We're getting there
- Announced in Glasgow'25 in May
- ...the we found some new problematic cases
- that took another year to fix
- Written for modern CPU and will have GPU acceleration (for hard cases)
- Results are *very good*, with excellent parallel scaling
- Handle hard cases particularly well
- Open source when done
- Stay tuned!





3 years later...

- We're getting there
- Announced in Glasgow'25 in May
- ...the we found some new problematic cases
- that took another year to fix
- Written for modern CPU and will have GPU acceleration (for hard cases)
- Results are *very* good, with excellent parallel scaling
- Handle hard cases particularly well
- Open source when done
- Stay tuned!





3 years later...

- We're getting there
- Announced in Glasgow'25 in May
- ...the we found some new problematic cases
- that took another year to fix
- Written for modern CPU and will have GPU acceleration (for hard cases)
- Results are *very* good, with excellent parallel scaling
- Handle hard cases particularly well
- Open source when done
- Stay tuned!





3 years later...

- We're getting there
- Announced in Glasgow'25 in May
- ...the we found some new problematic cases
- that took another year to fix
- Written for modern CPU and will have GPU acceleration (for hard cases)
- Results are *very* good, with excellent parallel scaling
- Handle hard cases particularly well
- Open source when done
- Stay tuned!



Thank you • شكرا



جامعة الملك عبد الله
للعلوم والتقنية
King Abdullah University of
Science and Technology