

GPUs and the computational efficiency of Gaussian process based models

Colin Rundel

Duke University

April 15, 2015

- 1 Background
- 2 Migratory Bird Spatial Assignment Model
- 3 Speciated PM_{2.5} Modeling
- 4 GPUs and Low Rank Approximations

The problem with GPs ...

Unless you are lucky (or clever), Gaussian process models are difficult to scale to large problems. For a Gaussian process $\mathcal{N}(\mu, \Sigma)$:

The problem with GPs ...

Unless you are lucky (or clever), Gaussian process models are difficult to scale to large problems. For a Gaussian process $\mathcal{N}(\mu, \Sigma)$:

Want a sample?

The problem with GPs ...

Unless you are lucky (or clever), Gaussian process models are difficult to scale to large problems. For a Gaussian process $\mathcal{N}(\mu, \Sigma)$:

Want a sample?

$$\mu + \text{Chol}(\Sigma) \times \mathbf{Z} \text{ with } Z_i \sim \mathcal{N}(0, 1)$$

The problem with GPs ...

Unless you are lucky (or clever), Gaussian process models are difficult to scale to large problems. For a Gaussian process $\mathcal{N}(\mu, \Sigma)$:

Want a sample?

$$\mu + \boxed{\text{Chol}(\Sigma)} \times \mathbf{Z} \text{ with } Z_i \sim \mathcal{N}(0, 1) \quad \mathcal{O}(n^3)$$

The problem with GPs ...

Unless you are lucky (or clever), Gaussian process models are difficult to scale to large problems. For a Gaussian process $\mathcal{N}(\mu, \Sigma)$:

Want a sample?

$$\mu + \boxed{\text{Chol}(\Sigma)} \times \mathbf{Z} \text{ with } Z_i \sim \mathcal{N}(0, 1) \quad \mathcal{O}(n^3)$$

Evaluate the (log) likelihood?

The problem with GPs ...

Unless you are lucky (or clever), Gaussian process models are difficult to scale to large problems. For a Gaussian process $\mathcal{N}(\mu, \Sigma)$:

Want a sample?

$$\mu + \boxed{\text{Chol}(\Sigma)} \times \mathbf{Z} \text{ with } Z_i \sim \mathcal{N}(0, 1) \quad \mathcal{O}(n^3)$$

Evaluate the (log) likelihood?

$$-\frac{1}{2} \log |\Sigma| - \frac{1}{2} (\mathbf{x} - \mu)' \Sigma^{-1} (\mathbf{x} - \mu) - \frac{n}{2} \log 2\pi$$

The problem with GPs ...

Unless you are lucky (or clever), Gaussian process models are difficult to scale to large problems. For a Gaussian process $\mathcal{N}(\mu, \Sigma)$:

Want a sample?

$$\mu + \text{Chol}(\Sigma) \times \mathbf{Z} \text{ with } Z_i \sim \mathcal{N}(0, 1) \quad \mathcal{O}(n^3)$$

Evaluate the (log) likelihood?

$$-\frac{1}{2} \log |\Sigma| - \frac{1}{2} (\mathbf{x} - \mu)' \Sigma^{-1} (\mathbf{x} - \mu) - \frac{n}{2} \log 2\pi \quad \mathcal{O}(n^3)$$

The problem with GPs ...

Unless you are lucky (or clever), Gaussian process models are difficult to scale to large problems. For a Gaussian process $\mathcal{N}(\mu, \Sigma)$:

Want a sample?

$$\mu + \text{Chol}(\Sigma) \times \mathbf{Z} \text{ with } Z_i \sim \mathcal{N}(0, 1) \quad \mathcal{O}(n^3)$$

Evaluate the (log) likelihood?

$$-\frac{1}{2} \log |\Sigma| - \frac{1}{2} (\mathbf{x} - \mu)' \Sigma^{-1} (\mathbf{x} - \mu) - \frac{n}{2} \log 2\pi \quad \mathcal{O}(n^3)$$

Update covariance parameter?

The problem with GPs ...

Unless you are lucky (or clever), Gaussian process models are difficult to scale to large problems. For a Gaussian process $\mathcal{N}(\mu, \Sigma)$:

Want a sample?

$$\mu + \text{Chol}(\Sigma) \times \mathbf{Z} \text{ with } Z_i \sim \mathcal{N}(0, 1) \quad \mathcal{O}(n^3)$$

Evaluate the (log) likelihood?

$$-\frac{1}{2} \log |\Sigma| - \frac{1}{2} (\mathbf{x} - \mu)' \Sigma^{-1} (\mathbf{x} - \mu) - \frac{n}{2} \log 2\pi \quad \mathcal{O}(n^3)$$

Update covariance parameter?

$$\{\Sigma\}_{ij} = \sigma^2 \exp(-\{d\}_{ij}\phi)$$

The problem with GPs ...

Unless you are lucky (or clever), Gaussian process models are difficult to scale to large problems. For a Gaussian process $\mathcal{N}(\mu, \Sigma)$:

Want a sample?

$$\mu + \text{Chol}(\Sigma) \times \mathbf{Z} \text{ with } Z_i \sim \mathcal{N}(0, 1) \quad \mathcal{O}(n^3)$$

Evaluate the (log) likelihood?

$$-\frac{1}{2} \log |\Sigma| - \frac{1}{2} (\mathbf{x} - \mu)' \Sigma^{-1} (\mathbf{x} - \mu) - \frac{n}{2} \log 2\pi \quad \mathcal{O}(n^3)$$

Update covariance parameter?

$$\{\Sigma\}_{ij} = \sigma^2 \exp(-\{d\}_{ij}\phi) \quad \mathcal{O}(n^2)$$

A simple guide to computational complexity

Linear complexity?

A simple guide to computational complexity

Linear complexity? - Go for it

A simple guide to computational complexity

Linear complexity? - Go for it

Quadratic complexity?

A simple guide to computational complexity

Linear complexity? - Go for it

Quadratic complexity? - Pray

A simple guide to computational complexity

Linear complexity? - Go for it

Quadratic complexity? - Pray

Cubic complexity?

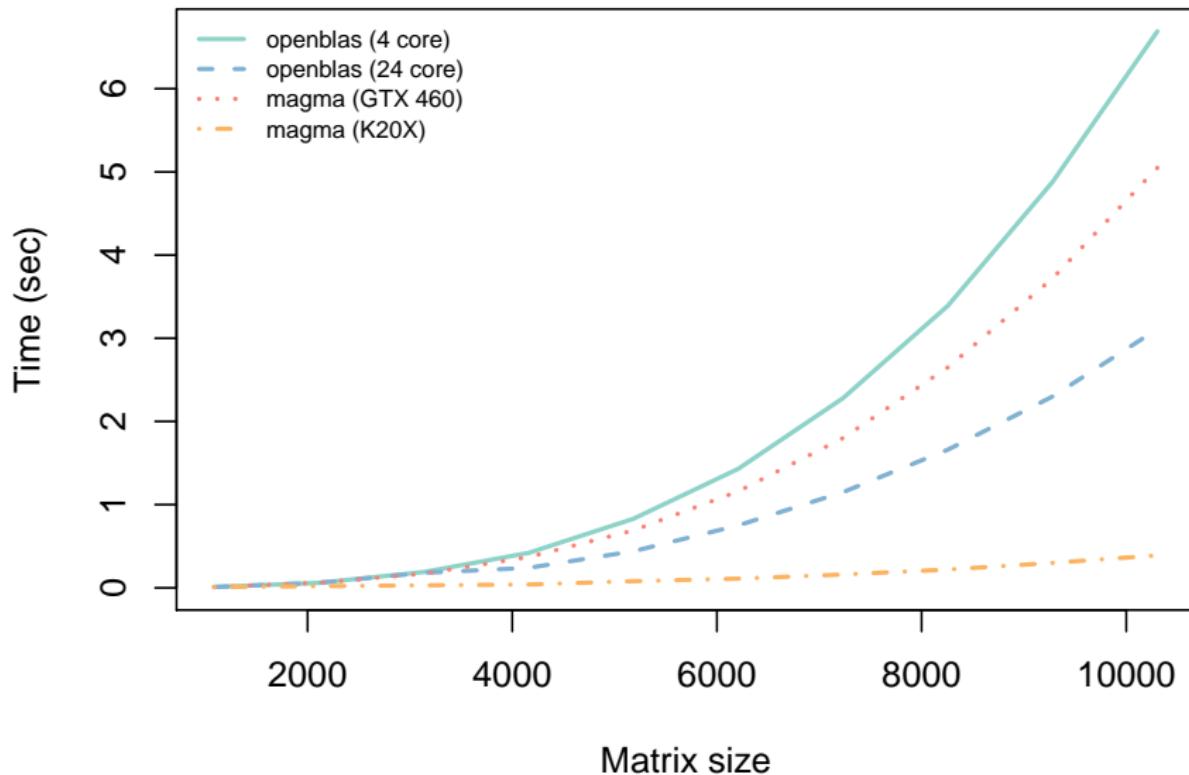
A simple guide to computational complexity

Linear complexity? - Go for it

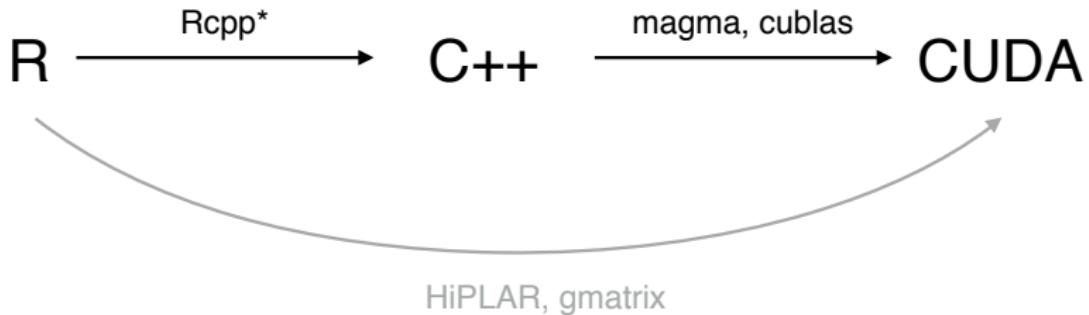
Quadratic complexity? - Pray

Cubic complexity? - Give up

Improving Cholesky



Tools and Optimization



Regardless of tools or workflow, measuring / profiling performance is critical.

- 1 Background
- 2 Migratory Bird Spatial Assignment Model
- 3 Speciated PM_{2.5} Modeling
- 4 GPUs and Low Rank Approximations

Background

Using intrinsic markers (genetic and isotopic signals) for the purpose of inferring migratory connectivity.

- Existing methods are too coarse for most applications
- Large amounts of data are available (>150,000 feather samples from >500 species)
- Genetic assignment methods are based on Wasser, et al. (2004)
- Isotopic assignment methods are based on Wunder, et al. (2005)

Data - DNA microsatellites and $\delta^2\text{H}$

Hermit Thrush (*Catharus guttatus*)

- 138 individuals
- 14 locations
- 6 loci
- 9-27 alleles / locus



Wilson's Warbler (*Wilsonia pusilla*)

- 163 individuals
- 8 locations
- 9 loci
- 15-31 alleles / locus



Allele Frequency Model

For the allele i , from locus l , at location k

$$\mathbf{y}_{\cdot lk} | \Theta \sim \text{Mult}(\sum_i y_{ilk}, \mathbf{f}_{lk})$$

$$f_{ilk} = \frac{\exp(\Theta_{ilk})}{\sum_i \exp(\Theta_{ilk})}$$

$$\Theta_{il} | \alpha, \mu \sim \mathcal{N}(\mu_{il}, \Sigma)$$

$$\{\Sigma\}_{ij} = \alpha_0 \exp\left(-(\{d\}_{ij}/\alpha_1)^{\alpha_2}\right) + \alpha_3 \mathbb{1}_{i=j}$$

Genetic Assignment Model

Assignment model using Hardy-Weinberg equilibrium allowing for genotyping (δ) and single amplification (γ) errors.

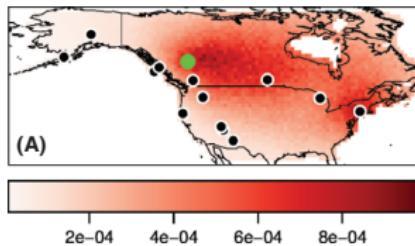
$$P(S_G | \mathbf{f}, k) = \prod_I P(i_I, j_I | \mathbf{f}, k)$$

$$P(i_I, j_I | \mathbf{f}, k) = \begin{cases} \gamma P(i_I | \mathbf{f}, k) + (1 - \gamma) P(i_I | \tilde{\mathbf{f}}, k)^2 & \text{if } i = j \\ (1 - \gamma) P(i_I | \mathbf{f}, k) P(j_I | \mathbf{f}, k) & \text{if } i \neq j \end{cases}$$

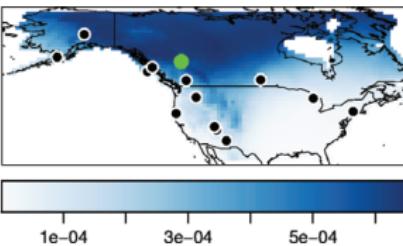
$$P(i_I | \mathbf{f}, k) = (1 - \delta) f_{lik} + \delta / m_I$$

Combined Model

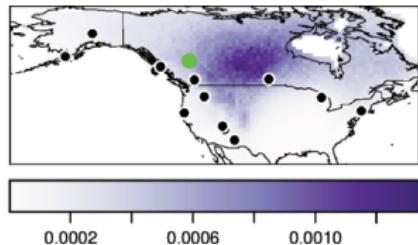
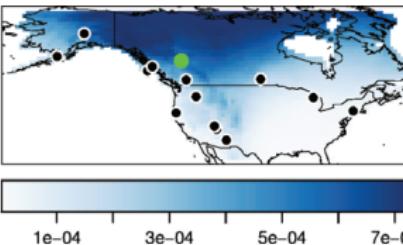
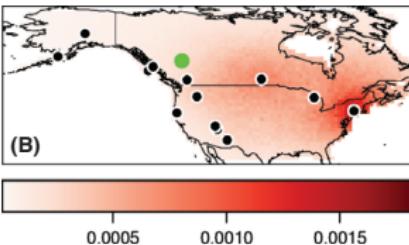
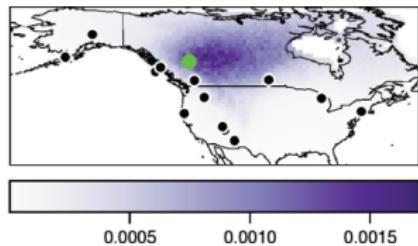
Genetic



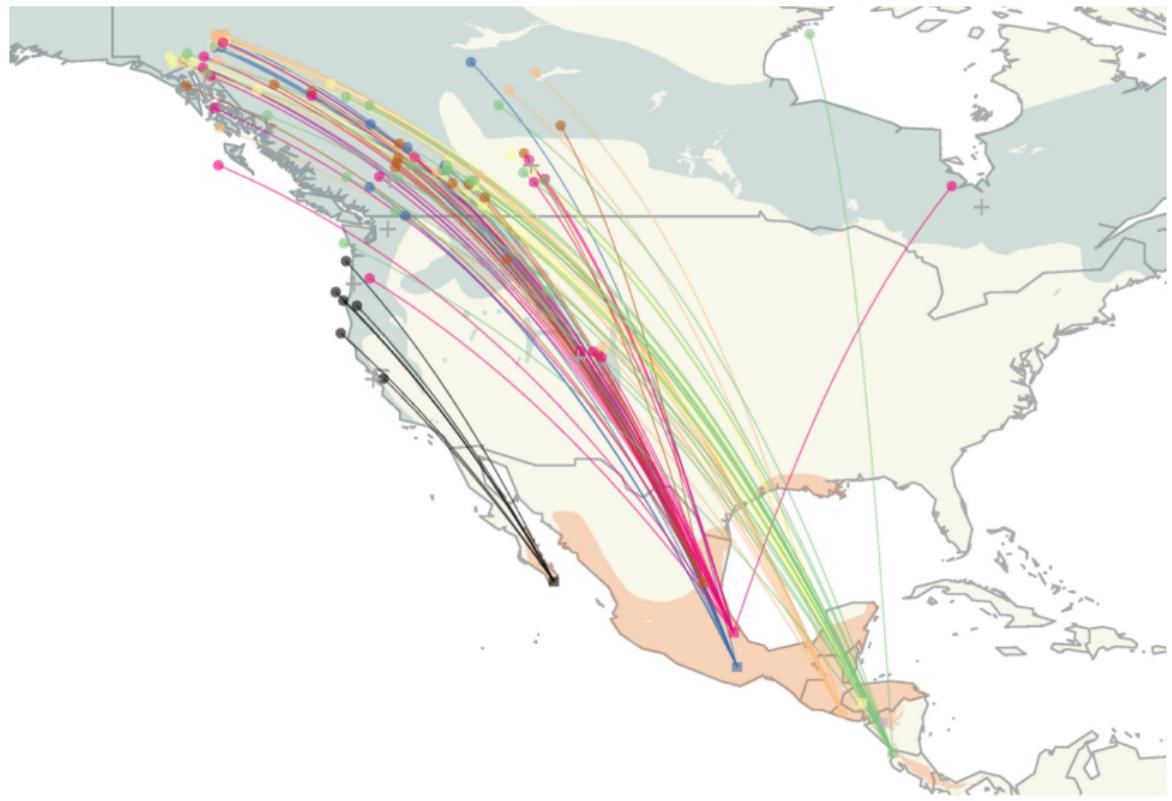
Isotopic



Combined



Migratory Connectivity



Implementation

Model fitting is done via MCMC (MH within Gibbs)

- Original implementation in pure C++ with minimal dependencies (Wasser, et al. (2004))
- Rewritten using R / C++ via Rcpp(Armadillo)
 - Code closer to matrix notation (and R)
 - Transparent use of high performance LAPACK implementations
 - R Package - isoscatR - <https://github.com/rundel/isoscatR>
- Model fitting performance is quite good
 - 300,000 iterations in ~ 5.5 minutes
- Bottleneck in drawing posterior predictive samples
 - 1,000 iterations in ~ 30 minutes

Prediction details

Why is the prediction slow?

Prediction details

Why is the prediction slow?

Predicting allele frequencies for Hermit thrush at 3318 novel locations.

Prediction details

Why is the prediction slow?

Predicting allele frequencies for Hermit thrush at 3318 novel locations.

To do so we sample from:

$$\Theta_p | \Theta_m \sim \mathcal{N}(\mu_p + \Sigma_{pm} \Sigma_m^{-1} (\Theta_m - \mu_m), \Sigma_p - \Sigma_{pm} \Sigma_m^{-1} \Sigma_{mp})$$

Prediction details

Why is the prediction slow?

Predicting allele frequencies for Hermit thrush at 3318 novel locations.

To do so we sample from:

$$\Theta_p | \Theta_m \sim \mathcal{N}(\mu_p + \Sigma_{pm}\Sigma_m^{-1}(\Theta_m - \mu_m), \Sigma_p - \Sigma_{pm}\Sigma_m^{-1}\Sigma_{mp})$$

Algorithm steps

- ① Calculate Σ_{pm} , Σ_p , and $\Sigma_p - \Sigma_{pm}\Sigma_m^{-1}\Sigma_{mp}$
- ② Calculate $\text{Chol}(\Sigma_p - \Sigma_{pm}\Sigma_m^{-1}\Sigma_{mp})$
- ③ Sample from MVN
- ④ Calculate allele frequencies

Posterior predictive sampling timings

Step	CPU (secs)	CPU+GPU (secs)	Rel. Performance
1. Covariances	1.080	0.046	23.0
2. Cholesky	0.467	0.208	2.3
3. Sample	0.049	0.052	0.9
4. Allele Freq	0.129	0.127	1.0
Total	1.732	0.465	3.7

Posterior predictive sampling timings

Step	CPU (secs)	CPU+GPU (secs)	Rel. Performance
1. Covariances	1.080	0.046	23.0
2. Cholesky	0.467	0.208	2.3
3. Sample	0.049	0.052	0.9
4. Allele Freq	0.129	0.127	1.0
Total	1.732	0.465	3.7

Total run time:

- CPU - 28.9 minutes
- CPU+GPU - 7.8 minutes

Posterior predictive sampling timings

Step	CPU (secs)	CPU+GPU (secs)	Rel. Performance
1. Covariances	1.080	0.046	23.0
2. Cholesky	0.467	0.208	2.3
3. Sample	0.049	0.052	0.9
4. Allele Freq	0.129	0.127	1.0
Total	1.732	0.465	3.7

Total run time:

- CPU - 28.9 minutes
- CPU+GPU - 7.8 minutes

\times CV runs $\begin{bmatrix} 166 \text{ for Hermit Thrush} \\ 179 \text{ for Wilson's Warbler} \end{bmatrix}$

Lessons

- Relatively small changes in one function resulted in 3 - 4x improvement
 - Cross validation results in two days instead of a week
 - 1-2 weeks of implementation, 1 week of tweaking / testing
 - Started with Cholesky, other optimizations followed

Lessons

- Relatively small changes in one function resulted in 3 - 4x improvement
 - Cross validation results in two days instead of a week
 - 1-2 weeks of implementation, 1 week of tweaking / testing
 - Started with Cholesky, other optimizations followed
- Issues
 - External library dependency makes package development (much) more complicated
 - Additional code verbosity and complexity

Improving Covariance Calculations

Covariance is assumed to be stationary and isotropic

- Elements of the covariance matrix can be calculated independently
- Small scale “embarrassingly parallel”
- Implementation is straight forward
(if we don't worry about things like symmetry)

```
__global__ void powered_exponential_cov_kernel(double* dist, double* cov,
                                                const int nm, const double sigma2,
                                                const double phi, const double kappa)
{
    int n_threads = gridDim.x * blockDim.x;
    int pos = blockDim.x * blockIdx.x + threadIdx.x;

    for (int i = pos; i < nm; i += n_threads)
    {
        cov[i] += sigma2 * exp( -pow(dist[i] * phi, kappa) );
    }
}
```

Building core tools

Common set of (expensive) tasks for GP models

- Covariance calculation
- Cholesky of Cov.
- Inverse of Cov.

Goal is to make performing these tasks on a GPU as painless as possible and allow interoperability with GPU (magma, CUBLAS) and CPU (Armadillo) libraries.

- GPU matrix class
- Modern resource management (RAII, move semantics)
- Simple translation between GPU and CPU memory

R Package - RcppGP - <https://github.com/rundel/RcppGP>

CPU vs GPU code

```
arma::mat prop_Sigma = arma::exp(-prop_phi * d_CIF);
arma::mat prop_Sigma_U = arma::chol(prop_Sigma);

double prop_Sigma_log_det = 2*arma::accu(arma::log(prop_Sigma_U.diag()));

arma::mat prop_Sigma_U_inv = arma::inv(arma::trimatu(prop_Sigma_U));
arma::mat prop_Sigma_inv = prop_Sigma_U_inv * prop_Sigma_U_inv.t();
```

```
exponential_cov_gpu(d_CIF_gpu.mat, cov_gpu.mat, nr_CIF, nr_CIF, 1.0, prop_phi, 64);
cov_gpu.chol('L',false);

double prop_Sigma_log_det = 2*arma::accu(arma::log(cov_gpu.get_mat().diag()));

cov_gpu.inv_chol('L',true);
arma::mat prop_Sigma_inv = cov_gpu.get_mat();
```

Back

- 1 Background
- 2 Migratory Bird Spatial Assignment Model
- 3 Speciated PM_{2.5} Modeling
- 4 GPUs and Low Rank Approximations

Background

Fine particulate matter ($\text{PM}_{2.5}$) is an EPA regulated air pollutant linked to a variety of adverse health effects

- Classified based on particle size ($< 2.5 \mu\text{m}$ diameter)
- Major species: Sulfate, Nitrate, Ammonium, Soil, Carbon.
- Minor species: trace elements (K, Mg, Ca), heavy metals (Cu, Fe), etc.
- Complex spatio-temporal dependence between species

Data

Speciated PM_{2.5} Sources

- Chemical Speciation Network (CSN) - 221 stations
- Interagency Monitoring of Protected Visual Environments (IMPROVE) - 172 stations

Total PM_{2.5} Sources

- Federal Reference Method (FRM) - 949 stations

Model Output

- Community Multi-scale Air Quality (CMAQ) - 12 km grid

Data Issues

- Monitoring frequency
- Total vs Sum of Species

Species Model Details

For the 5 major species (Sulfate, Nitrate, Ammonium, Soil, Carbon) and the two networks (CSN, IMPROVE):

$$C_t^i(\mathbf{s}) = Z_t^i(\mathbf{s}) + \epsilon_{C,t}^i(\mathbf{s})$$

$$I_t^i(\mathbf{s}) = Z_t^i(\mathbf{s}) + \epsilon_{I,t}^i(\mathbf{s})$$

where $Z_t^i(\mathbf{s})$ are the latent “true” concentrations of species i at time t and locations \mathbf{s} , and is given by

$$Z_t^i(\mathbf{s}) = \max(0, \tilde{Z}_t^i(\mathbf{s}))$$

$$\tilde{Z}_t^i(\mathbf{s}) = \beta_{0,t}^i + \beta_{0,t}^i(\mathbf{s}) + \beta_{1,t}^i Q_t^i(B_s)$$

Total PM_{2.5} Model Details

For total PM_{2.5} from the three networks (CSN, IMPROVE, FRM):

$$C_t^{tot}(\mathbf{s}) = Z_t^{tot}(\mathbf{s}) + \epsilon_{C,t}^{tot}(\mathbf{s})$$

$$I_t^{tot}(\mathbf{s}) = Z_t^{tot}(\mathbf{s}) + \epsilon_{I,t}^{tot}(\mathbf{s})$$

$$F_t^{tot}(\mathbf{s}) = Z_t^{tot}(\mathbf{s}) + \epsilon_{F,t}^{tot}(\mathbf{s})$$

where $Z_t^{tot}(\mathbf{s})$ are the latent “true” concentration of total PM_{2.5} at time t and locations \mathbf{s} , which is given by the sum of the major species and the “other” species concentrations.

$$Z_t^{tot}(\mathbf{s}) = \sum_{i=1}^5 Z_t^i(\mathbf{s}) + Z_t^o(\mathbf{s})$$

$$Z_t^o(s) = \max \left(0, \tilde{Z}_t^o(s) \right) \quad \tilde{Z}_t^o(s) = \beta_{0,t}^o + \beta_{0,t}^o(\mathbf{s}) + \beta_{1,t}^o Q_t^o(B_s)$$

Spatial Dependence

Spatial dependence enters the model through the $\beta_{0,t}^i(s)$ parameters for $i \in \{o, 1, 2, 3, 4, 5\}$.

$$\beta_{0,t}^i(\mathbf{s}) = \sigma_t^i w_t^i(\mathbf{s})$$

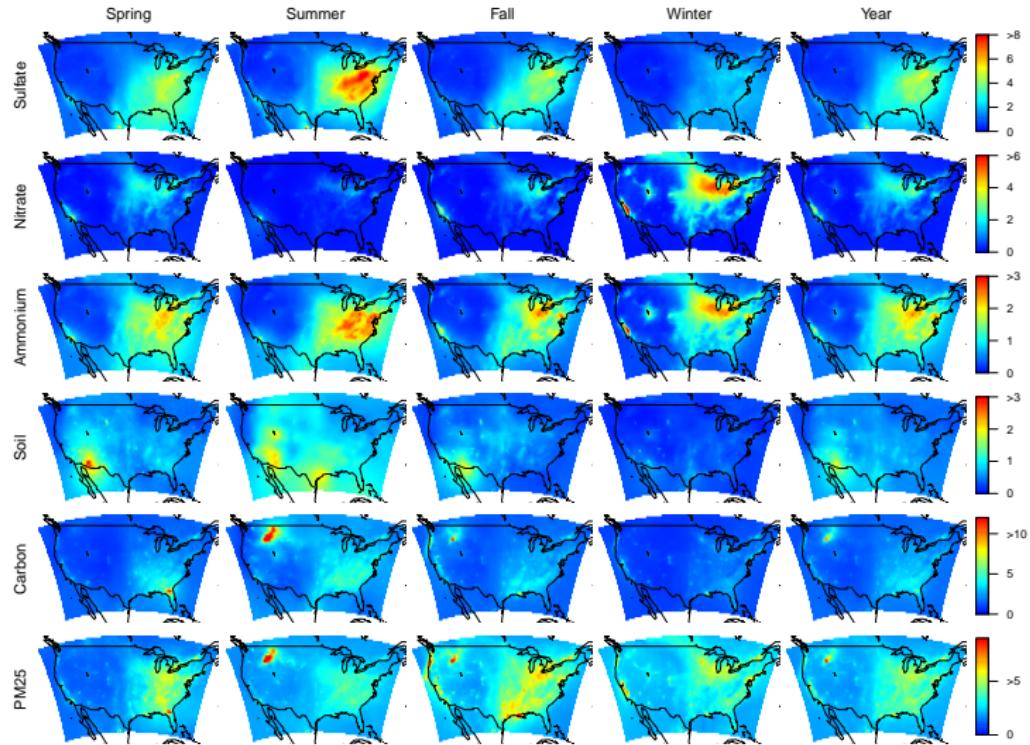
where $w_t^i(\mathbf{s})$ are zero mean, variance 1, Gaussian processes with exponential correlation given by

$$\text{corr}(w_t^i(\mathbf{s}), w_t^i(\mathbf{s}')) = \exp(-\phi_t^i |\mathbf{s} - \mathbf{s}'|)$$

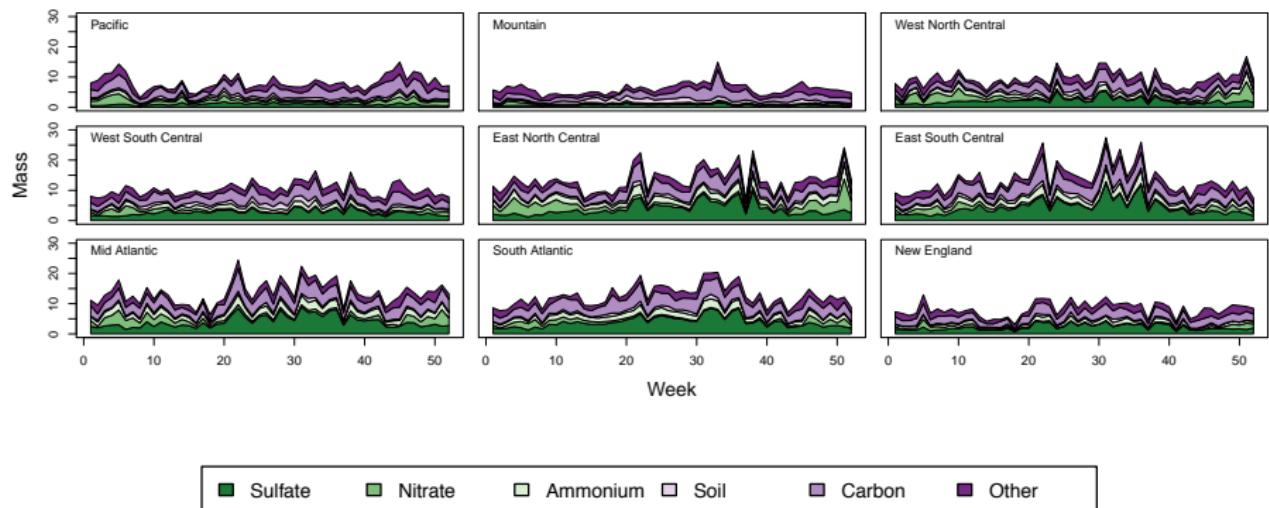
Additional dependence between species is introduced via coregionalization,

$$\begin{pmatrix} \beta_{0,t}^i(\mathbf{s}) \\ \beta_{0,t}^j(\mathbf{s}) \end{pmatrix} = \mathbf{A}_t \begin{pmatrix} w_t^i(\mathbf{s}) \\ w_t^j(\mathbf{s}) \end{pmatrix}.$$

Model results



Model results



MCMC performance

Parameter	CPU (secs)	CPU+GPU (secs)	Rel. Performance
β_0, β_1	0.00029	0.00030	0.97
$\beta_0(s)$	0.09205	0.09132	1.00
σ^2	0.00383	0.00385	0.99
ϕ	0.46084	0.25174	1.83
τ_i^2, τ_{tot}^2	0.00003	0.00003	1.00
Total	0.55708	0.34729	1.60

Run times

Total run time for model fitting (50,000 iterations):

- CPU - 7.7 hours \times 52 weeks
- CPU+GPU - 4.8 hours

Run times

Total run time for model fitting (50,000 iterations):

- CPU - 7.7 hours × 52 weeks
- CPU+GPU - 4.8 hours

Total run time for model prediction at 5950 locations (1,000 iterations):

- CPU - 7.2 hours × 52 weeks
- CPU+GPU - 4.3 hours

Run times

Total run time for model fitting (50,000 iterations):

- CPU - 7.7 hours × 52 weeks
- CPU+GPU - 4.8 hours

Total run time for model prediction at 5950 locations (1,000 iterations):

- CPU - 7.2 hours × 52 weeks
- CPU+GPU - 4.3 hours

One run takes about 775 hours (4.6 days) total on CPU alone, 473 (2.8 days) on CPU and GPU.

Run times

Total run time for model fitting (50,000 iterations):

- CPU - 7.7 hours × 52 weeks
- CPU+GPU - 4.8 hours

Total run time for model prediction at 5950 locations (1,000 iterations):

- CPU - 7.2 hours × 52 weeks
- CPU+GPU - 4.3 hours

One run takes about 775 hours (4.6 days) total on CPU alone, 473 (2.8 days) on CPU and GPU.

 × 3 model variants
 × 10 for cross validation

Lessons

- Established infrastructure makes a huge difference in development time
 - 1 hour to go from CPU implementation to CPU+GPU implementation
 - Code shown previously is 2/3 of the changes necessary [Code](#)
- In practice, was easier to run CPU only code across more servers (configuration time / effort)
 - Not possible (or at least easy) for models variants that are not independent in time.
 - There will be ~ 20 desktops with GPUs available in the department (available via Condor)
- Rcpp Attributes offer huge advantages in development and deployment
 - Simplifies external dependencies (locally)
 - Full compatibility is the goal for RcppGP

- 1 Background
- 2 Migratory Bird Spatial Assignment Model
- 3 Speciated PM_{2.5} Modeling
- 4 GPUs and Low Rank Approximations

Low rank approximations

For a Gaussian process

$$Y(\mathbf{s}) = \mathbf{x}(\mathbf{s})' \boldsymbol{\beta} + w(\mathbf{s}) + \epsilon, \quad \epsilon \sim N(0, \tau^2 I)$$

$$w(\mathbf{s}) \sim N(0, \mathbf{C}(\mathbf{s})), \quad \mathbf{C}(\mathbf{s}, \mathbf{s}') = \sigma \rho(\mathbf{s}, \mathbf{s}' | \theta)$$

we can approximate $\mathbf{C}(\mathbf{s})$ with a low rank approximation with the form $\mathbf{U} \mathbf{S} \mathbf{V}'$ where \mathbf{U} and \mathbf{V} are $n \times k$ and \mathbf{S} is $k \times k$.

Low rank approximations

For a Gaussian process

$$Y(\mathbf{s}) = \mathbf{x}(\mathbf{s})' \boldsymbol{\beta} + w(\mathbf{s}) + \epsilon, \quad \epsilon \sim N(0, \tau^2 I)$$

$$w(\mathbf{s}) \sim N(0, \mathbf{C}(\mathbf{s})), \quad \mathbf{C}(\mathbf{s}, \mathbf{s}') = \sigma \rho(\mathbf{s}, \mathbf{s}' | \theta)$$

we can approximate $\mathbf{C}(\mathbf{s})$ with a low rank approximation with the form $\mathbf{U} \mathbf{S} \mathbf{V}'$ where \mathbf{U} and \mathbf{V} are $n \times k$ and \mathbf{S} is $k \times k$.

This allows for the use of the Sherman-Morrison-Woodbury formula for the inverse (and determinant),

$$\mathbf{C}(\mathbf{s})^{-1} \approx (\mathbf{A} + \mathbf{U} \mathbf{S} \mathbf{V}')^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{U} (\mathbf{S}^{-1} + \mathbf{V}' \mathbf{A}^{-1} \mathbf{U})^{-1} \mathbf{V}' \mathbf{A}^{-1}.$$

Gaussian Predictive Processes

For a rank k approximation,

- Pick k knot locations \mathbf{s}^*
- Calculate knot covariance ($\mathbf{C}(\mathbf{s}^*)$) and knot cross-covariance ($(\mathbf{C}(\mathbf{s}^*))^{-1}$)
- Approximate full covariance

$$\mathbf{C}(\mathbf{s}) \approx \mathbf{C}(\mathbf{s}, \mathbf{s}^*) \mathbf{C}(\mathbf{s}^*)^{-1} \mathbf{C}(\mathbf{s}^*, \mathbf{s}).$$

- Systematically underestimates variance, inflates τ^2 .
- Modified predictive process corrects this using

$$\mathbf{C}(\mathbf{s}) \approx \mathbf{C}(\mathbf{s}, \mathbf{s}^*) \mathbf{C}(\mathbf{s}^*)^{-1} \mathbf{C}(\mathbf{s}^*, \mathbf{s}) + \text{diag}\left(\mathbf{C}(\mathbf{s}) - \mathbf{C}(\mathbf{s}, \mathbf{s}^*) \mathbf{C}(\mathbf{s}^*)^{-1} \mathbf{C}(\mathbf{s}^*, \mathbf{s})\right).$$

Banerjee, Gelfand, Finley, Sang (2008) Finley, Sang, Banerjee, Gelfand (2008)

Low Rank Approximations via Random Projections

- ① Starting with an $m \times n$ matrix \mathbf{A} .
- ② Draw an $n \times k + p$ Gaussian random matrix Ω .
- ③ Form $\mathbf{Y} = \mathbf{A}\Omega$ and compute its QR factorization $\mathbf{Y} = \mathbf{Q} \mathbf{R}$
- ④ Form the $k + p \times n$ matrix $\mathbf{B} = \mathbf{Q}' \mathbf{A}$.
- ⑤ Compute the SVD of the small matrix \mathbf{B} , $\mathbf{B} = \hat{\mathbf{U}} \mathbf{S} \mathbf{V}'$.
- ⑥ Form the matrix $\mathbf{U} = \mathbf{Q} \hat{\mathbf{U}}$.

Low Rank Approximations via Random Projections

- ① Starting with an $m \times n$ matrix \mathbf{A} .
- ② Draw an $n \times k + p$ Gaussian random matrix Ω .
- ③ Form $\mathbf{Y} = \mathbf{A}\Omega$ and compute its QR factorization $\mathbf{Y} = \mathbf{Q} \mathbf{R}$
- ④ Form the $k + p \times n$ matrix $\mathbf{B} = \mathbf{Q}' \mathbf{A}$.
- ⑤ Compute the SVD of the small matrix \mathbf{B} , $\mathbf{B} = \hat{\mathbf{U}} \mathbf{S} \mathbf{V}'$.
- ⑥ Form the matrix $\mathbf{U} = \mathbf{Q} \hat{\mathbf{U}}$.

Resulting approximation has nicely bounded expected error,

$$\mathbb{E} \ \| \mathbf{A} - \mathbf{USV}' \| \leq \left[1 + \frac{4\sqrt{k+p}}{p-1} \sqrt{\min(m,n)} \right] \sigma_{k+1}.$$

Halko, Martinsson, Tropp (2011)

Random Matrix Low Rank Decompositions and GPs

Preceeding algorithm can be modified slightly to take advantage of the positive definite structure of a covariance matrix.

- ① Starting with an $n \times n$ covariance matrix \mathbf{A} .
- ② Draw an $n \times k + p$ Gaussian random matrix Ω .
- ③ Form $\mathbf{Y} = \mathbf{A}\Omega$ and compute its QR factorization $\mathbf{Y} = \mathbf{Q}\mathbf{R}$
- ④ Form the $k + p \times k + p$ matrix $\mathbf{B} = \mathbf{Q}'\mathbf{A}\mathbf{Q}$.
- ⑤ Compute the eigen decomposition of the small matrix \mathbf{B} ,
$$\mathbf{B} = \hat{\mathbf{U}}\mathbf{S}\hat{\mathbf{U}}'$$
.
- ⑥ Form the matrix $\mathbf{U} = \mathbf{Q}\hat{\mathbf{U}}$.

Once again we have a bound on the error,

$$\mathbb{E} \|\mathbf{A} - \mathbf{Q}(\mathbf{Q}'\mathbf{A}\mathbf{Q})\mathbf{Q}'\| = \mathbb{E} \|\mathbf{A} - \mathbf{U}\mathbf{S}\mathbf{U}'\| \lesssim c \cdot \sigma_{k+1}.$$

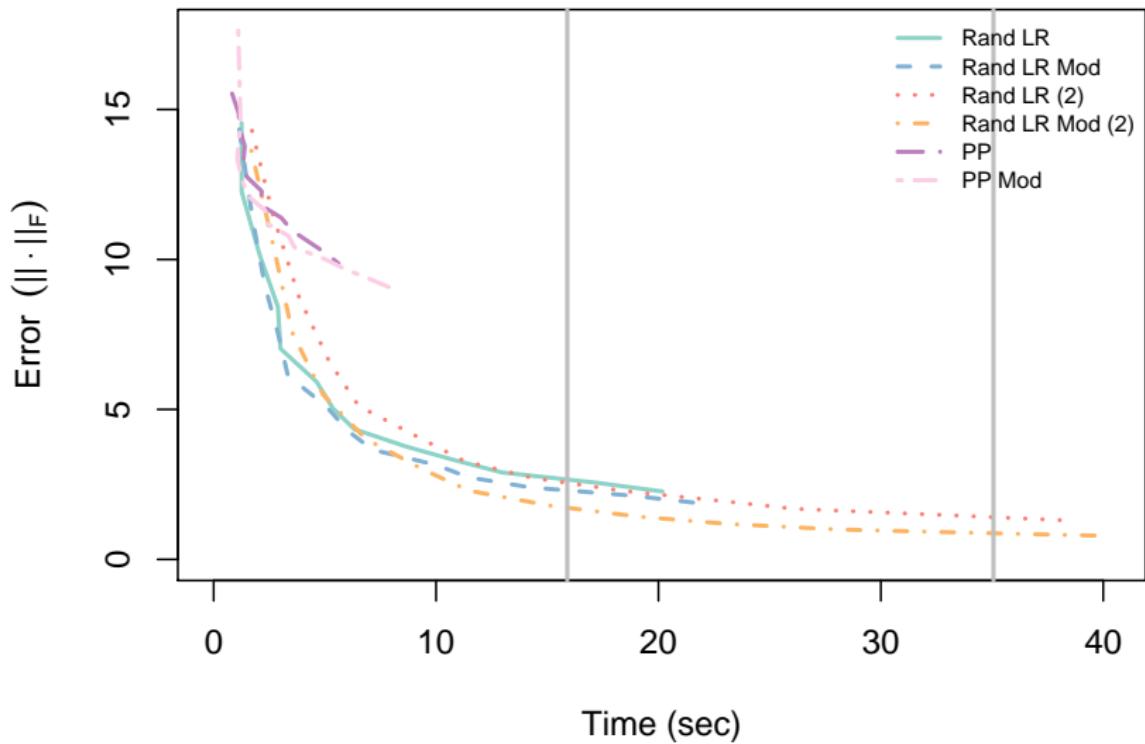
Halko, Martinsson, Tropp (2011), Banerjee, Dunson, Tokdar (2012)

Low Rank Approximations and GPUs

Both predictive process and random matrix low rank approximations are good candidates for acceleration using GPUs.

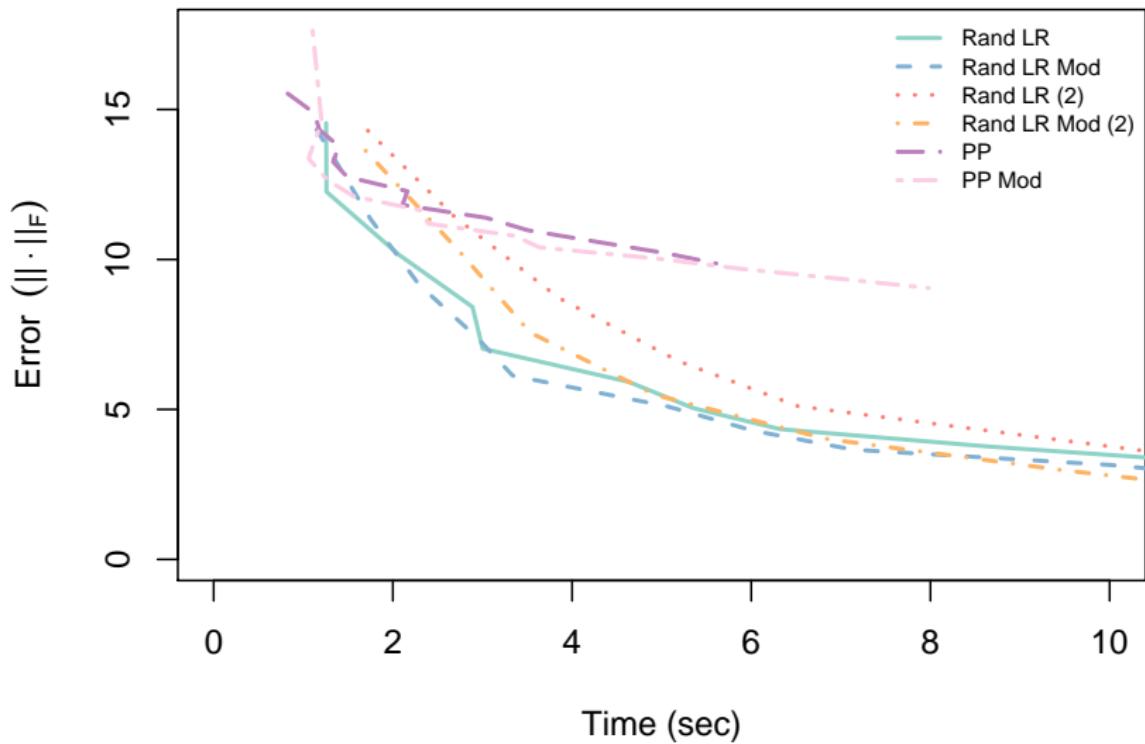
- Both use Sherman-Woodbury-Morrison to calculate the inverse (involves matrix multiplication, addition, and a small matrix inverse).
- Predictive processes involves several covariance matrix calculations (knots and cross-covariance) and a small matrix inverse.
- Random matrix low rank involves a large matrix multiplication ($\mathbf{A}\Omega$) and several small matrix decompositions (QR, eigen).
- Functionality for both approaches included in current version of RcppGP (inv_lr and inv_pp).

Matrix inverse (fixed rank, strong dependence)



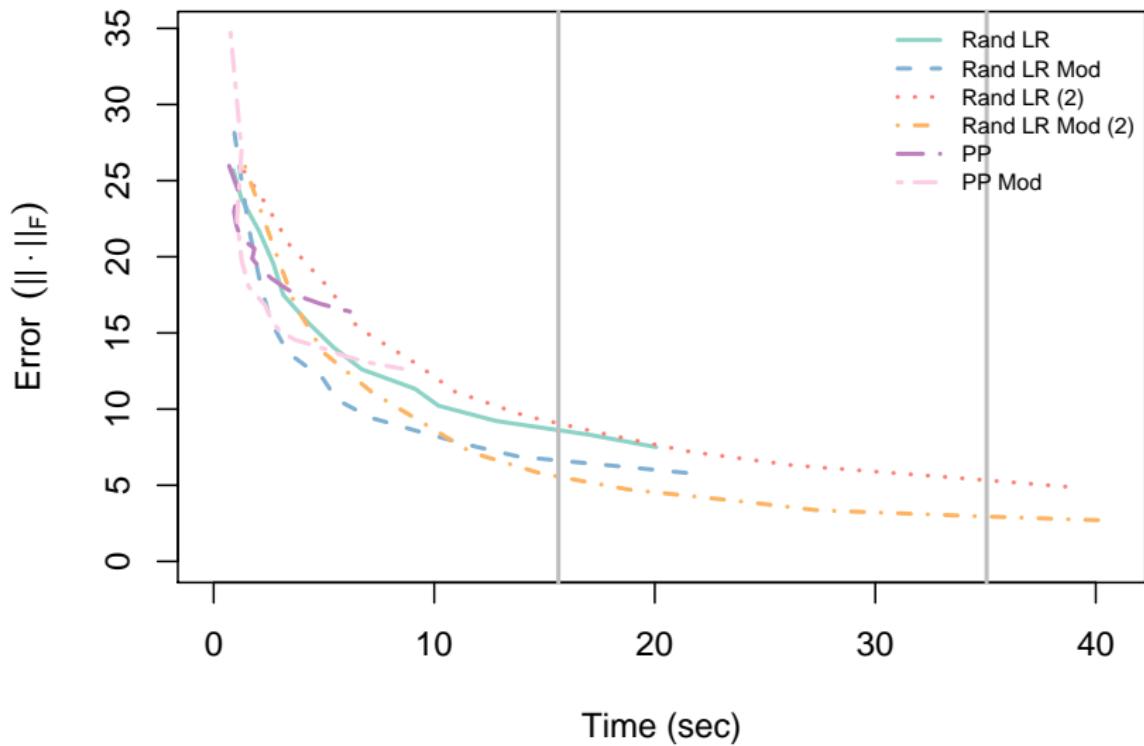
$$n = 15000, \quad k = \{100, \dots, 4900\}$$

Matrix inverse (fixed rank, strong dependence)



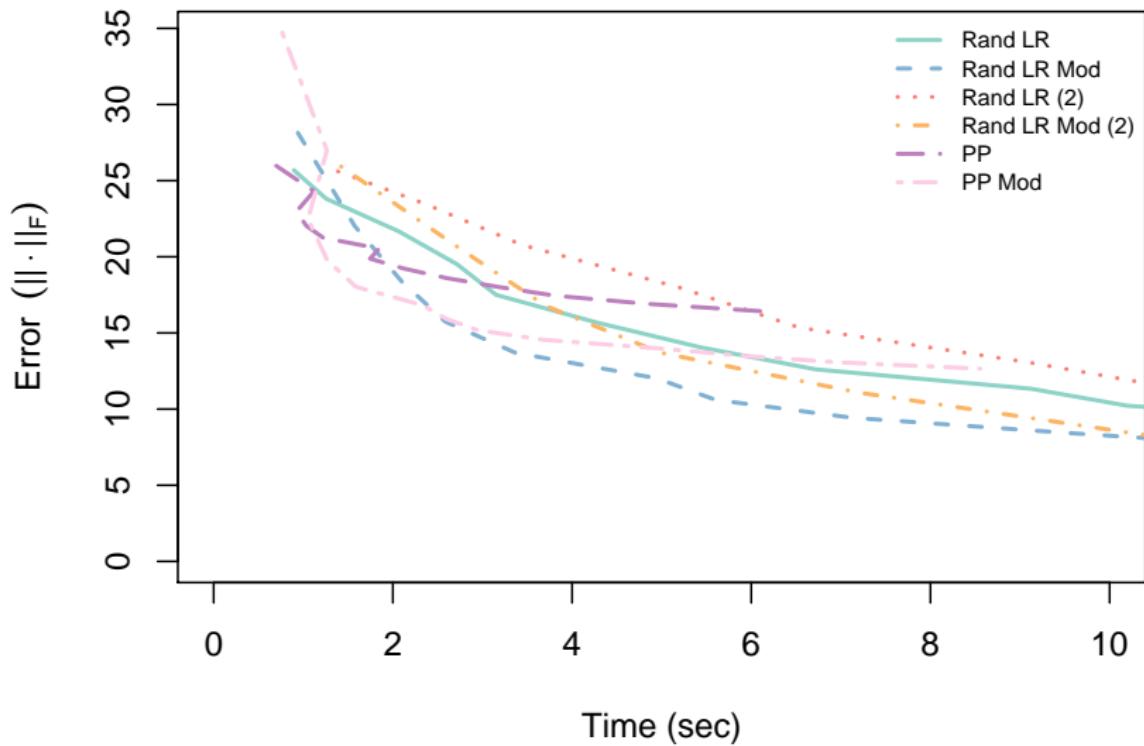
$$n = 15000, \quad k = \{100, \dots, 4900\}$$

Matrix inverse (fixed rank, weak dependence)



$$n = 15000, \quad k = \{100, \dots, 4900\}$$

Matrix inverse (fixed rank, weak dependence)



$$n = 15000, \quad k = \{100, \dots, 4900\}$$

Rand. Matrix Low Rank Decompositions for Prediction

This approach can also be used for prediction, if we want to sample

$$\mathcal{N}(0, \Sigma) \text{ with } \Sigma \approx \mathbf{U} \mathbf{S} \mathbf{U}' = (\mathbf{U} \mathbf{S}^{1/2} \mathbf{U}') (\mathbf{U} \mathbf{S}^{1/2} \mathbf{U}')'$$

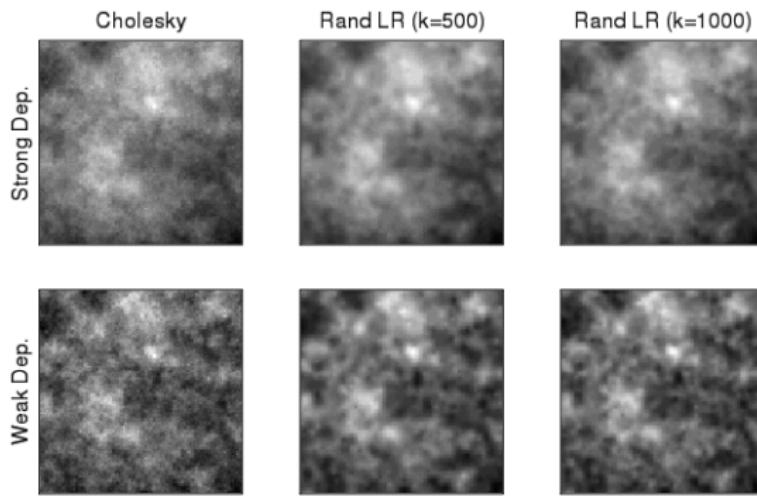
then $X_{pred} = (\mathbf{U} \mathbf{S}^{1/2} \mathbf{U}') \times \mathbf{Z}$ where $Z_i \sim \mathcal{N}(0, 1)$.

Rand. Matrix Low Rank Decompositions for Prediction

This approach can also be used for prediction, if we want to sample

$$\mathcal{N}(0, \Sigma) \text{ with } \Sigma \approx \mathbf{U} \mathbf{S} \mathbf{U}' = (\mathbf{U} \mathbf{S}^{1/2} \mathbf{U}') (\mathbf{U} \mathbf{S}^{1/2} \mathbf{U}')'$$

then $X_{pred} = (\mathbf{U} \mathbf{S}^{1/2} \mathbf{U}') \times \mathbf{Z}$ where $Z_i \sim \mathcal{N}(0, 1)$.



$$n = 1000, \quad p = 10000$$

Dehdari, Deutsch (2012)

Future Directions

- Refinement of RcppGP
 - Transition to header only implementation
 - Transparent GPU to CPU failover
 - Support for fixed error (instead of rank) random matrix low rank decomposition
 - Thinking about out-of-memory based approaches
- Future of GPUs, CUDA, and Magma
 - Single vs. Multi-GPU algorithms
 - Mixed precision algorithms
 - NVBLAS
 - Unified memory
 - cuSolver

Acknowledgments

Migratory Connectivity

- John Novembre - UChicago
- Thomas Smith - UCLA
- Kristen Ruegg - UCLA, UCSC
- Center for Tropical Research,
UCLA IoES

Speciated PM_{2.5}

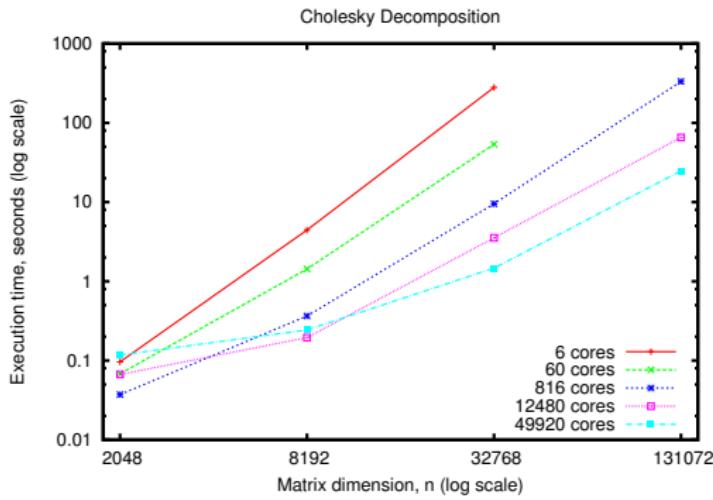
- Alan Gelfand - Duke
- Dave Holland - EPA
- Erin Schliep - Duke

RcppGP <https://github.com/rundel/RcppGP>
Talk <https://github.com/rundel/Presentations>

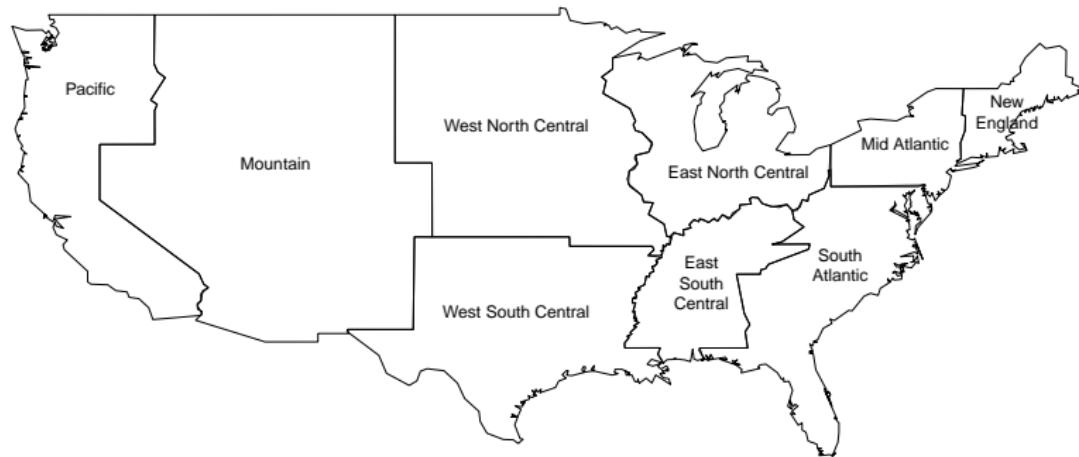
Another Approach

bigGP is an R package written by Chris Paciorek, et al.

- Specialized implementation of LA operation for GPs
- Designed to run on large super computer clusters
- Uses both shared and distributed memory
- Able to fit models on the order of $n = 65k$ (32 GB Cov. matrix)



Regions



Back