## Lecture 22

Computational Methods for GPs

Colin Rundel
04/12/2017

# GPs and Computational Complexity

## The problem with GPs

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

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

Want to sample $\boldsymbol{y}$?

$$\boldsymbol{\mu} + \boxed{\text{Chol}(\boldsymbol{\Sigma})} \times \boldsymbol{z} \text{ with } z_i \sim \mathcal{N}(0, 1) \qquad \mathcal{O}\left(n^3\right)$$

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

Want to sample $y$?

$$\boldsymbol{\mu} + \boxed{\text{Chol}(\boldsymbol{\Sigma})} \times z \text{ with } z_i \sim \mathcal{N}(0, 1) \qquad \mathcal{O}\left(n^3\right)$$

Evaluate the (log) likelihood?

$$-\frac{1}{2} \log \boxed{|\boldsymbol{\Sigma}|} - \frac{1}{2}(x - \boldsymbol{\mu})' \boxed{\boldsymbol{\Sigma}^{-1}} (x - \boldsymbol{\mu}) - \frac{n}{2} \log 2\pi \qquad \mathcal{O}\left(n^3\right)$$

# The problem with GPs

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

Want to sample $\boldsymbol{y}$?

$$\boldsymbol{\mu} + \boxed{\text{Chol}(\boldsymbol{\Sigma})} \times \boldsymbol{z} \text{ with } z_i \sim \mathcal{N}(0,1) \qquad \mathcal{O}\left(n^3\right)$$

Evaluate the (log) likelihood?

$$-\frac{1}{2}\log \boxed{|\Sigma|} - \frac{1}{2}(x - \boldsymbol{\mu})' \boxed{\boldsymbol{\Sigma}^{-1}}(x - \boldsymbol{\mu}) - \frac{n}{2}\log 2\pi \qquad \mathcal{O}\left(n^3\right)$$

Update covariance parameter?

$$\boxed{\{\Sigma\}_{ij}} = \sigma^2 \exp(-\{d\}_{ij}\phi) + \sigma_n^2 \, 1_{i=j} \qquad \mathcal{O}\left(n^2\right)$$

$\mathcal{O}(n)$ - Linear complexity

$\mathcal{O}\left(n\right)$ - Linear complexity - Go for it

$\mathcal{O}\left(n\right)$ - Linear complexity - Go for it

$\mathcal{O}\left(n^2\right)$ - Quadratic complexity

$\mathcal{O}(n)$ - Linear complexity - Go for it

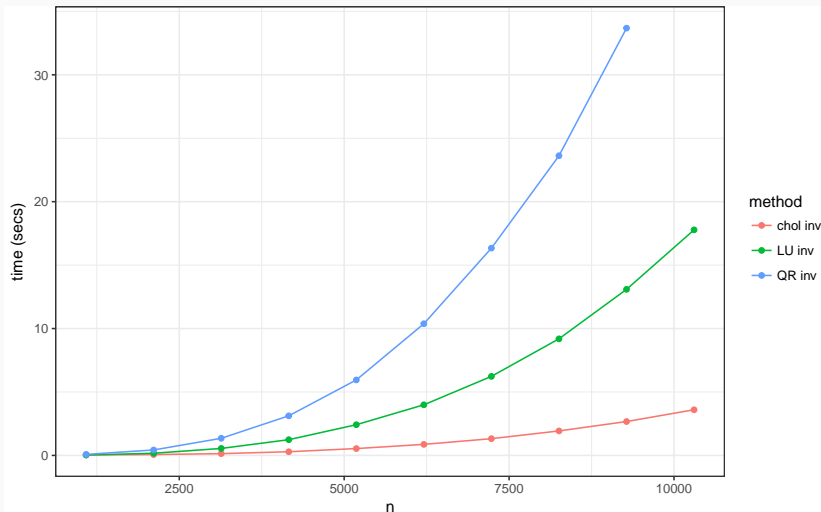$\mathcal{O}(n^2)$ - Quadratic complexity - Pray

$\mathcal{O}\left(n\right)$ - Linear complexity - Go for it

$\mathcal{O}\left(n^2\right)$ - Quadratic complexity - Pray

$\mathcal{O}\left(n^3\right)$ - Cubic complexity

$\mathcal{O}\left(n\right)$ - Linear complexity - Go for it

$\mathcal{O}\left(n^2\right)$ - Quadratic complexity - Pray

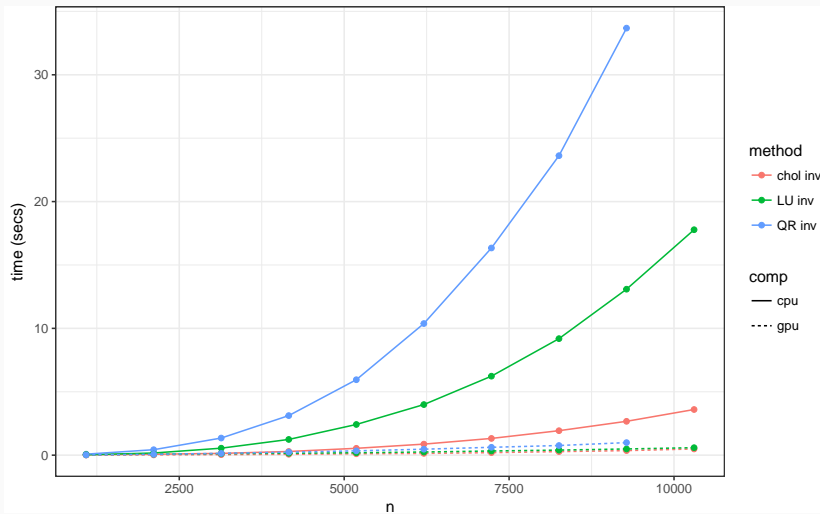$\mathcal{O}\left(n^3\right)$ - Cubic complexity - Give up

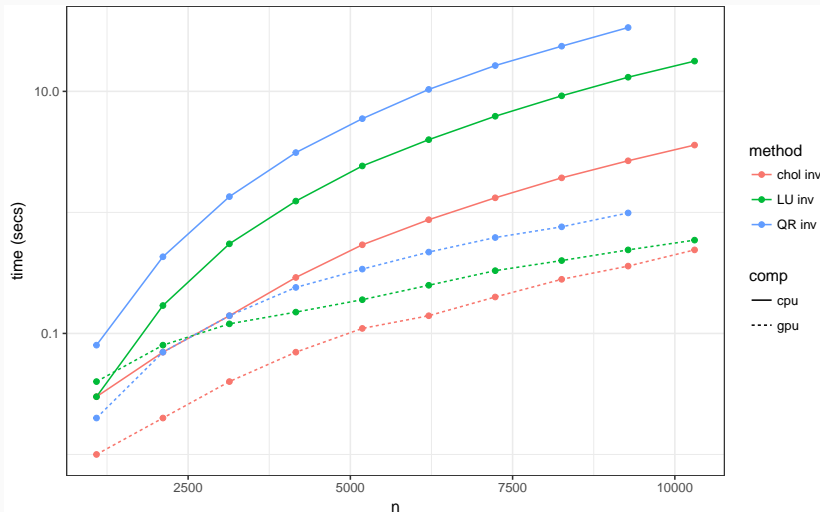After fitting the GP need to sample from the posterior predictive distribution at $\sim 3000$ locations

$$y_p \sim \mathcal{N}\left(\mu_p + \Sigma_{po}\Sigma_o^{-1}(y_o - \mu_o),\ \Sigma_p - \Sigma_{po}\Sigma_o^{-1}\Sigma_{op}\right)$$

After fitting the GP need to sample from the posterior predictive distribution at $\sim$ 3000 locations

$$y_p \sim \mathcal{N}\left(\mu_p + \Sigma_{po}\Sigma_o^{-1}(y_o - \mu_o),\ \Sigma_p - \Sigma_{po}\Sigma_o^{-1}\Sigma_{op}\right)$$

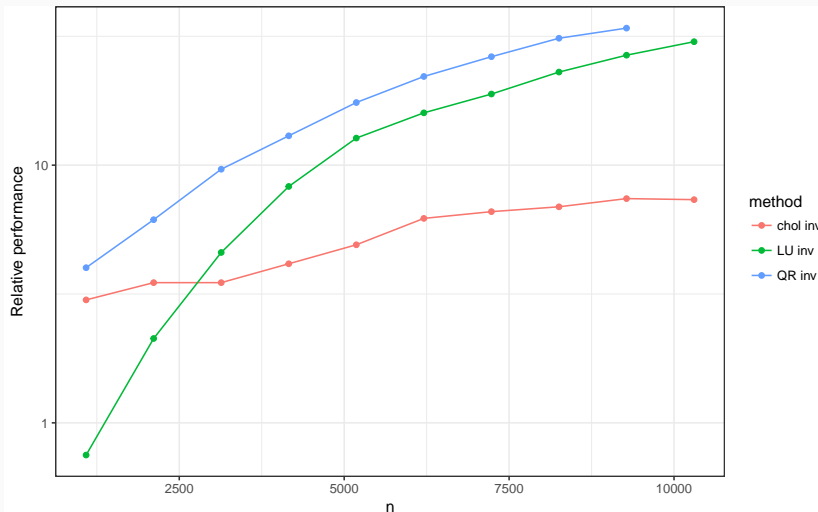|     | Step | CPU (secs) | CPU+GPU (secs) | Rel. Performance |
|-----|------|-----------|----------------|------------------|
| 1.  | Calc. $\Sigma_p$, $\Sigma_{po}$ | 1.080 | 0.046 | 23.0 |
| 2.  | Calc. chol($\Sigma_p - \Sigma_{po}\Sigma_o^{-1}\Sigma_{op}$) | 0.467 | 0.208 | 2.3 |
| 3.  | Calc. $\mu_{p\|o}$ + chol($\Sigma_{p\|o}$) $\times Z$ | 0.049 | 0.052 | 0.9 |
| 4.  | Calc. Allele Prob | 0.129 | 0.127 | 1.0 |
|     | Total | 1.732 | 0.465 | 3.7 |

Total run time: CPU (28.9 min), CPU+GPU (7.8 min)

# Cholesky CPU vs GPU (P100)
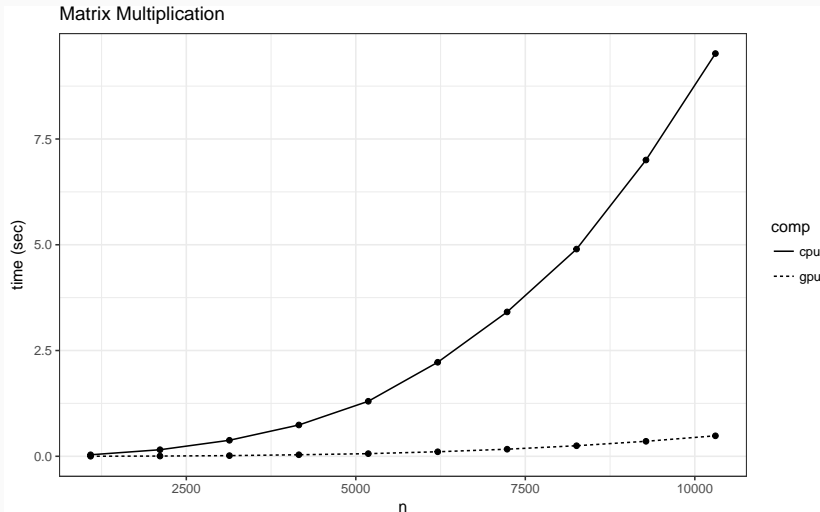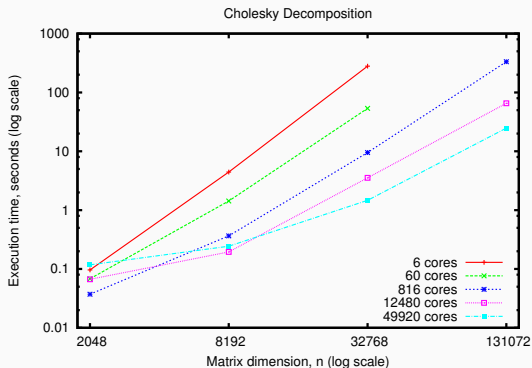
Matrix Multiplication

## An even bigger hammer

`bigGP` is an R package written by Chris Paciorek (UC Berkeley), et al.

- Specialized distributed implementation of linear algebra 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 $n = 65$k (32 GB Cov. matrix)

- Spectral domain / basis functions

- Covariance tapering

- GMRF approximations

- Low-rank approximations

- Nearest-neighbor models

# Low Rank Approximations

## Low rank approximations in general

Lets look at the example of the singular value decomposition of a matrix,

$$\underset{n\times m}{M} \;=\; \underset{n\times n}{U}\ \underset{n\times m}{\mathrm{diag}(S)}\ \underset{m\times m}{V^{t}}$$

where $U$ are called the left singular vectors, $V$ the right singular vectors, and $S$ the singular values. Usually the singular values and vectors are ordered such that the signular values are in descending order.

## Low rank approximations in general

Lets look at the example of the singular value decomposition of a matrix,

$$\underset{n \times m}{M} = \underset{n \times n}{U} \underset{n \times m}{\operatorname{diag}(S)} \underset{m \times m}{V^t}$$

where *U* are called the left singular vectors, *V* the right singular vectors, and *S* the singular values. Usually the singular values and vectors are ordered such that the signular values are in descending order.

It turns out (Eckart–Young theorem) that we can approximate *M* as having rank *r* by defining $\tilde{S}$ to only have the *r* largest singular values (others set to zero).

$$\underset{n \times m}{\tilde{M}} = \underset{n \times n}{U} \underset{n \times m}{\operatorname{diag}(\tilde{S})} \underset{m \times m}{V^t} = \underset{n \times k}{\tilde{U}} \underset{k \times k}{\operatorname{diag}(\tilde{S})} \underset{k \times m}{\tilde{V}^t}$$

# Example

$$M = \begin{pmatrix} 1.000 & 0.500 & 0.333 & 0.250 \\ 0.500 & 0.333 & 0.250 & 0.200 \\ 0.333 & 0.250 & 0.200 & 0.167 \\ 0.250 & 0.200 & 0.167 & 0.143 \end{pmatrix} = U \operatorname{diag}(S) V^t$$

$$U = V = \begin{pmatrix} -0.79 & 0.58 & -0.18 & -0.03 \\ -0.45 & -0.37 & 0.74 & 0.33 \\ -0.32 & -0.51 & -0.10 & -0.79 \\ -0.25 & -0.51 & -0.64 & 0.51 \end{pmatrix}$$

$$S = \begin{pmatrix} 1.50 & 0.17 & 0.01 & 0.00 \end{pmatrix}$$

## Example

$$M = \begin{pmatrix} 1.000 & 0.500 & 0.333 & 0.250 \\ 0.500 & 0.333 & 0.250 & 0.200 \\ 0.333 & 0.250 & 0.200 & 0.167 \\ 0.250 & 0.200 & 0.167 & 0.143 \end{pmatrix} = U \operatorname{diag}(S) \, V^t$$

$$U = V = \begin{pmatrix} -0.79 & 0.58 & -0.18 & -0.03 \\ -0.45 & -0.37 & 0.74 & 0.33 \\ -0.32 & -0.51 & -0.10 & -0.79 \\ -0.25 & -0.51 & -0.64 & 0.51 \end{pmatrix}$$

$$S = \begin{pmatrix} 1.50 & 0.17 & 0.01 & 0.00 \end{pmatrix}$$

Rank 2 approximation:

$$\tilde{M} = \begin{pmatrix} -0.79 & 0.58 \\ -0.45 & -0.37 \\ -0.32 & -0.51 \\ -0.25 & -0.51 \end{pmatrix} \begin{pmatrix} 1.50 & 0.00 \\ 0.00 & 0.17 \end{pmatrix} \begin{pmatrix} -0.79 & -0.45 & -0.32 & -0.25 \\ 0.58 & -0.37 & -0.51 & -0.51 \end{pmatrix}$$

$$= \begin{pmatrix} 1.000 & 0.501 & 0.333 & 0.249 \\ 0.501 & 0.330 & 0.251 & 0.203 \\ 0.333 & 0.251 & 0.200 & 0.166 \\ 0.249 & 0.203 & 0.166 & 0.140 \end{pmatrix}$$
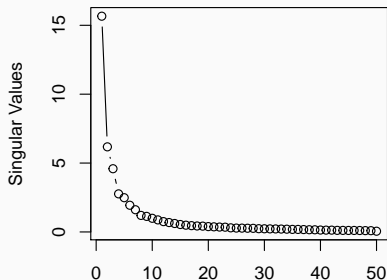
## Approximation Error

We can measure the error of the approximation using the Frobenius norm,

$$\|M - \tilde{M}\|_F = \left( \sum_{i=1}^{m} \sum_{j=1}^{n} (M_{ij} - \tilde{M}_{ij})^2 \right)^{1/2}$$

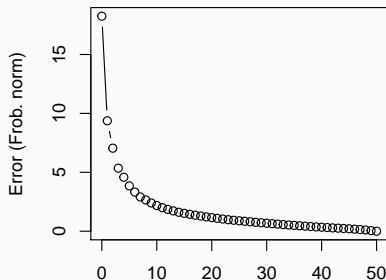We can measure the error of the approximation using the Frobenius norm,

$$\|M - \tilde{M}\|_F = \left( \sum_{i=1}^{m} \sum_{j=1}^{n} (M_{ij} - \tilde{M}_{ij})^2 \right)^{1/2}$$

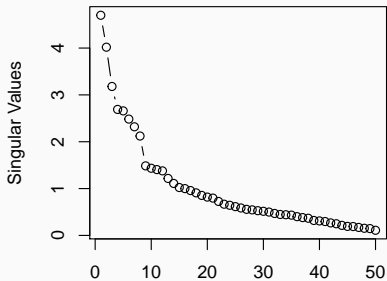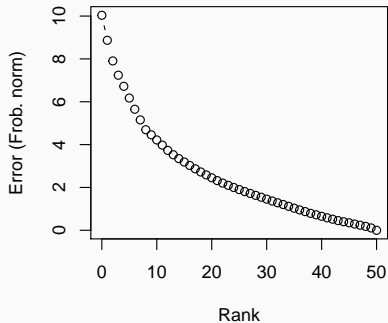Strong dependence (large eff. range):

Weak dependence (short eff. range):

## How does this help? (Sherman-Morrison-Woodbury)

There is an immensely useful linear algebra identity, the
Sherman-Morrison-*Woodbury* formula, for the inverse (and determinant) of
a decomposed matrix,

$$
\underset{n \times m}{\tilde{M}}{}^{-1} = \left( \underset{n \times m}{A} + \underset{n \times k}{U} \underset{k \times k}{S} \underset{k \times m}{V^t} \right)^{-1}
$$
$$
= A^{-1} - A^{-1}U \left( S^{-1} + V^t A^{-1} U \right)^{-1} V^t A^{-1}.
$$

## How does this help? (Sherman-Morrison-Woodbury)

There is an immensely useful linear algebra identity, the Sherman-Morrison-*Woodbury* formula, for the inverse (and determinant) of a decomposed matrix,

$$\underset{n\times m}{\tilde{M}}{}^{-1} = \left( \underset{n\times m}{A} + \underset{n\times k}{U}\, \underset{k\times k}{S}\, \underset{k\times m}{V^t} \right)^{-1}$$
$$= A^{-1} - A^{-1}U \left( S^{-1} + V^t A^{-1} U \right)^{-1} V^t A^{-1}.$$

How does this help?

- Imagine that $A = \text{diag}(A)$, then it is trivial to find $A^{-1}$.
- $S^{-1}$ is $k \times k$ which is hopefully small, or even better $S = \text{diag}(S)$.
- $\left( S^{-1} + V^t A^{-1} U \right)$ is $k \times k$ which is hopefully small.

## Aside - Determinant

Remember for any MVN distribution when evaluating the likelihood

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

we need the inverse of $\Sigma$ as well as its *determinant*.

Remember for any MVN distribution when evaluating the likelihood

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

we need the inverse of $\Sigma$ as well as its *determinant*.

- For a full rank Cholesky decomposition we get the determinant for "free".

$$|M| = |LL^t| = \prod_{i=1}^{n} \left(\text{diag}(L)_i\right)^2$$

## Aside - Determinant

Remember for any MVN distribution when evaluating the likelihood

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

we need the inverse of $\Sigma$ as well as its *determinant*.

· For a full rank Cholesky decomposition we get the determinant for "free".

$$|M| = |LL^t| = \prod_{i=1}^{n}\left(\text{diag}(L)_i\right)^2$$

· For a low rank approximation the Sherman-Morrison-Woodbury / Determinant lemma gives us,

$$
\begin{aligned}
\det(\tilde{M}) &= \det(A + USV^t) \\
&= \det(S^{-1} + V^tA^{-1}U) \ \det(S) \ \det(A)
\end{aligned}
$$

19

## Low rank approximations for GPs

For a standard spatial random effects model,

$$y(s) = x(s)\,\beta + w(s) + \epsilon, \quad \epsilon \sim N(0,\,\tau^2 I)$$
$$w(s) \sim \mathcal{N}(0,\,\Sigma(s)), \quad \Sigma(s, s') = \sigma\,\rho(s, s'|\theta)$$

if we can replace $\Sigma(s)$ with a low rank approximation of the form

- $\Sigma(s) \approx U\,S\,V^t$ where
- $U$ and $V$ are $n \times k$,
- $S$ is $k \times k$, and
- $A = \tau^2 I$ or a similar diagonal matrix

Predictive Processes

## Gaussian Predictive Processes

For a rank $k$ approximation,

- Pick $k$ knot locations $s^\star$
- Calculate knot covariance, $\Sigma(s^\star)$, and knot cross-covariance, $\Sigma(s, s^\star)$
- Approximate full covariance using

$$\Sigma(s) \approx \underset{n \times k}{\Sigma(s, s^\star)} \, \underset{k \times k}{\Sigma(s^\star)^{-1}} \, \underset{k \times n}{\Sigma(s^\star, s)}.$$

- PPs systematically underestimates variance ($\sigma^2$) and inflate $\tau^2$, Modified predictive processs corrects this using

$$\begin{aligned}
\Sigma(s) \approx &\Sigma(s, s^\star) \, \Sigma(s^\star)^{-1} \, \Sigma(s^\star, s) \\
&+ \text{diag}\Big( \Sigma(s) - \Sigma(s, s^\star) \, \Sigma(s^\star)^{-1} \, \Sigma(s^\star, s) \Big).
\end{aligned}$$

Banerjee, Gelfand, Finley, Sang (2008)    Finley, Sang, Banerjee, Gelfand (2008)
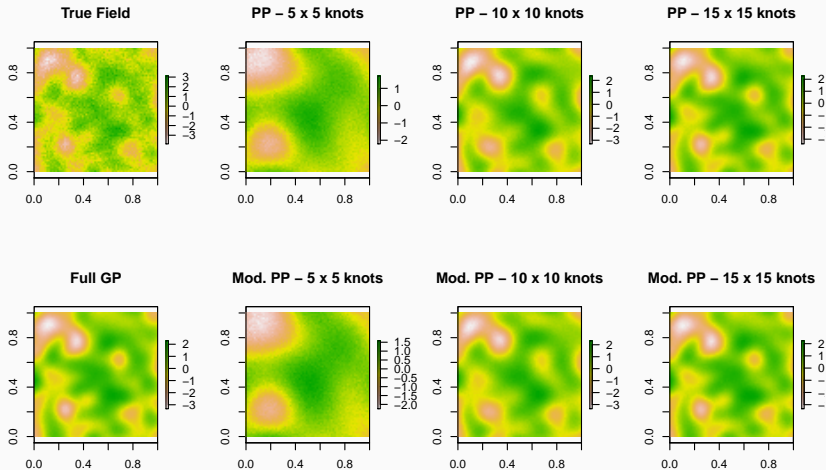
Below we have a surface generate from a squared exponential Gaussian
Process where

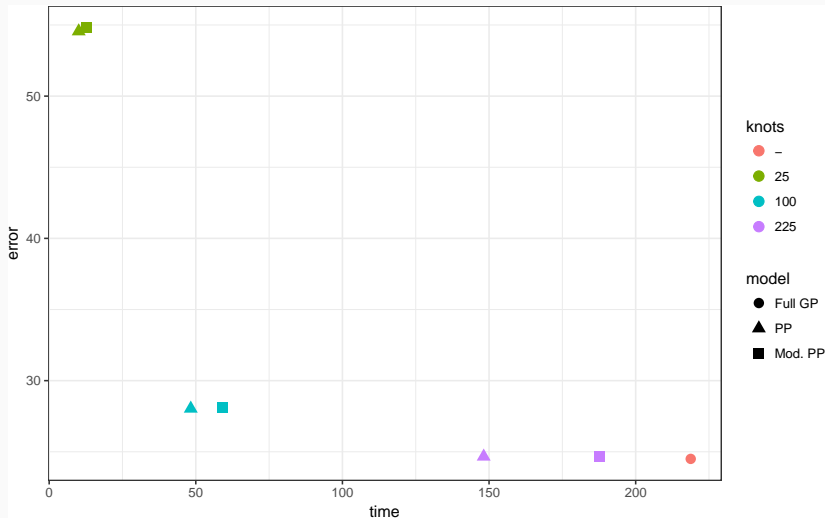$$\{\Sigma\}_{ij} = \sigma^2 \exp\left(-(\phi\, d)^2\right) + \tau^2 I$$
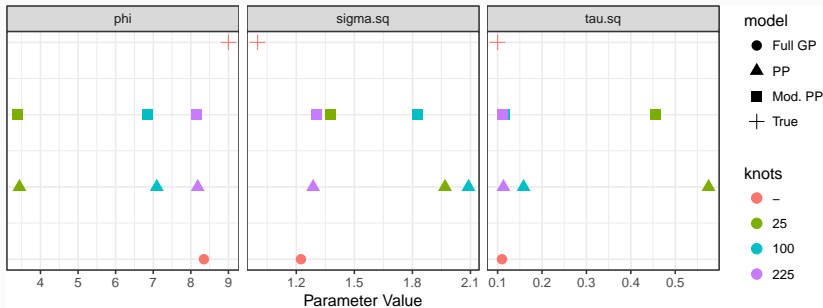
$$\sigma^2 = 1 \quad \phi = 9 \quad \tau^2 = 0.1$$

# Random Projections

## Low Rank Approximations via Random Projections

1. Starting with an $m \times n$ matrix $A$.
2. Draw an $n \times k + p$ Gaussian random matrix $\Omega$.
3. Form $Y = A\Omega$ and compute its QR factorization $Y = QR$

4. Form the $k + p \times n$ matrix $B = Q'A$.
5. Compute the SVD of the small matrix $B$, $B = \hat{U}SV'$.
6. Form the matrix $U = Q\hat{U}$.

Resulting approximation has a bounded expected error,

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

Halko, Martinsson, Tropp (2011)

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

Once again we have a bound on the error,

$$\mathsf{E}\,\|A - Q(Q'AQ)Q'\| = \mathsf{E}\,\|A - USU'\| \lesssim c \cdot \sigma_{k+1}.$$
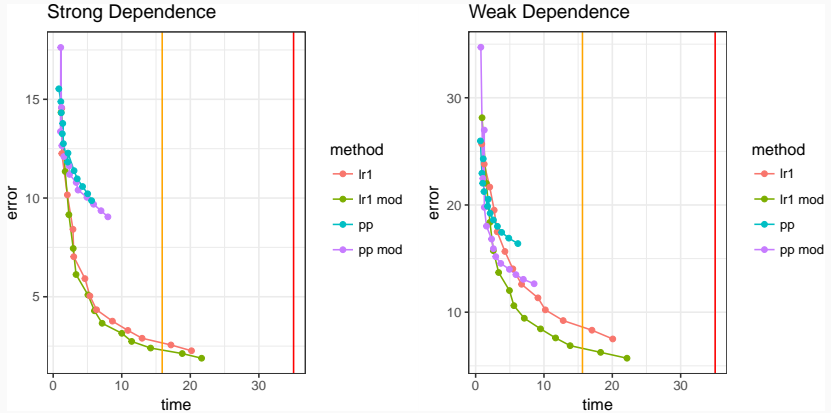
Halko, Martinsson, Tropp (2011), Banerjee, Dunson, Tokdar (2012)

## 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 approximations involves a large matrix multiplication ($A\,\Omega$) and several small matrix decompositions (QR, eigen).

30

Strong Dependence

Weak Dependence

## Rand. Projection LR Decompositions for Prediction

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

$$y \sim \mathcal{N}(0, \boldsymbol{\Sigma})$$
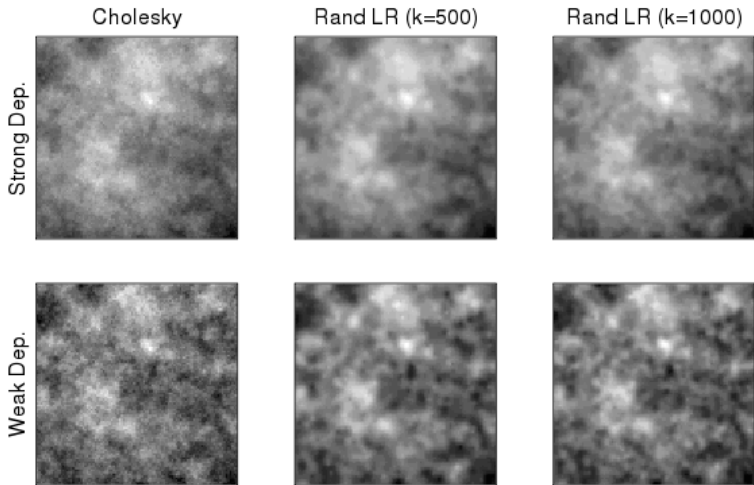$$\Sigma \approx USU^t = (US^{1/2}U^t)(US^{1/2}U^t)^t$$

then

$$y_{\text{pred}} = (U S^{1/2} U^t) \times Z \text{ where } Z_i \sim \mathcal{N}(0, 1)$$

because $U^t U = I$ since $U$ is an orthogonal matrix.

Dehdari, Deutsch (2012)

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

33