

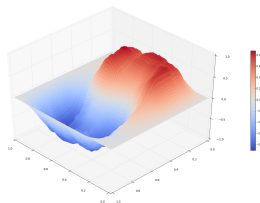
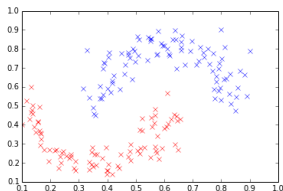
GPU-Based Regression Analysis on Sparse Grids

Steffen Hirschmann

September 24, 2014

- ▶ Feature space Ω d -dimensional
- ▶ Assign a label $k \in \mathbb{R}$ to $\vec{x} \in \Omega$
- ▶ Supervised learning. Pre-classified data set:

$$S = \{(\vec{x}_i, y_i) \in \Omega \times \mathbb{R}\}_{i=1}^m$$

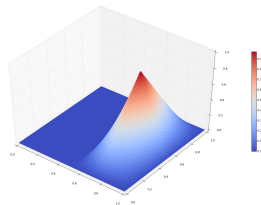
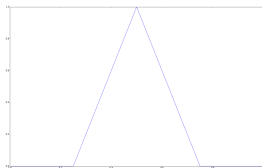


Approach on Sparse Grids

- ▶ Discretize Ω using sparse grids
- ▶ Piecewise d -linear sparse grid function

$$u(\vec{x}) = \sum_{i=1}^g v_i \phi_i(\vec{x})$$

- ▶ $\phi_i(\vec{x})$ tensor products of 1d hat functions.



- ▶ Determine u via a least squares problem

$$u = \operatorname{argmin}_u \left(\frac{1}{m} \sum_{i=1}^m (y_i - u(\vec{x}_i))^2 + \lambda \sum_{i=1}^g v_i^2 \right)$$

- ▶ Minimization by setting $\frac{\partial}{\partial v_i}$ to zero.
- ▶ Resulting SLE:

$$\left(\frac{1}{m} BB^T + \lambda I \right) \vec{v} = \frac{1}{m} B\vec{y}$$

- ▶ $B_{ij} = \phi_i(\vec{x}_j)$
- ▶ Solve via CG; handled by framework SG⁺⁺

$$\underbrace{\begin{pmatrix} \phi_1(\vec{x}_1) & \phi_1(\vec{x}_2) & \dots & \phi_1(\vec{x}_m) \\ \phi_2(\vec{x}_1) & \phi_2(\vec{x}_2) & \dots & \phi_2(\vec{x}_m) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_g(\vec{x}_1) & \phi_g(\vec{x}_2) & \dots & \phi_g(\vec{x}_m) \end{pmatrix}}_B \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_m \end{pmatrix} = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_g \end{pmatrix}$$
$$\underbrace{\begin{pmatrix} \phi_1(\vec{x}_1) & \phi_2(\vec{x}_1) & \dots & \phi_g(\vec{x}_1) \\ \phi_1(\vec{x}_2) & \phi_2(\vec{x}_2) & \dots & \phi_g(\vec{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(\vec{x}_m) & \phi_2(\vec{x}_m) & \dots & \phi_g(\vec{x}_m) \end{pmatrix}}_{B^T} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_g \end{pmatrix} = \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_m \end{pmatrix}$$

Implementation

Basic Algorithm - Parallelization

$$\underbrace{\begin{pmatrix} \phi_1(\vec{x}_1) & \phi_1(\vec{x}_2) & \dots & \phi_1(\vec{x}_m) \\ \phi_2(\vec{x}_1) & \phi_2(\vec{x}_2) & \dots & \phi_2(\vec{x}_m) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_g(\vec{x}_1) & \phi_g(\vec{x}_2) & \dots & \phi_g(\vec{x}_m) \end{pmatrix}}_B \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_m \end{pmatrix} = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_g \end{pmatrix}$$
$$\underbrace{\begin{pmatrix} \phi_1(\vec{x}_1) & \phi_2(\vec{x}_1) & \dots & \phi_g(\vec{x}_1) \\ \phi_1(\vec{x}_2) & \phi_2(\vec{x}_2) & \dots & \phi_g(\vec{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(\vec{x}_m) & \phi_2(\vec{x}_m) & \dots & \phi_g(\vec{x}_m) \end{pmatrix}}_{B^T} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_g \end{pmatrix} = \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_m \end{pmatrix}$$

Workload of one thread (index t):

Operation B :

$$\langle B_{t,\cdot}, \mathbf{v} \rangle = \sum_{j=1}^m v_j \phi_t(\vec{x}_j)$$

Operation B^T :

$$\langle B_{\cdot,t}, \mathbf{v} \rangle = \sum_{j=1}^g \alpha_j \phi_j(\vec{x}_t)$$

Workload of one thread (index t):

Operation B :

$$\langle B_{t,\cdot}, v \rangle = \sum_{j=1}^m v_j \prod_d \phi_{l_{t,d}, i_{t,d}}^{1D}(x_{j,d})$$

Operation B^T :

$$\langle B_{\cdot,t}, v \rangle = \sum_{j=1}^g \alpha_j \prod_d \phi_{l_{j,d}, i_{j,d}}^{1D}(x_{t,d})$$

Workload of one thread (index t):

Operation B :

$$\langle B_{t,\cdot}, v \rangle = \sum_{j=1}^m v_j \prod_d \phi_{l_{t,d}, i_{t,d}}^{1D}(x_{j,d})$$

Operation B^T :

$$\langle B_{\cdot,t}, v \rangle = \sum_{j=1}^g \alpha_j \prod_d \phi_{l_{j,d}, i_{j,d}}^{1D}(x_{t,d})$$

- ▶ Different number of addends (\rightarrow deferred to memory management)

Workload of one thread (index t):

Operation B :

$$\langle B_{t,\cdot}, v \rangle = \sum_{j=1}^m v_j \prod_d \phi_{t,d,i_{t,d}}^{1D}(x_{j,d})$$

Operation B^T :

$$\langle B_{\cdot,t}, v \rangle = \sum_{j=1}^g \alpha_j \prod_d \phi_{j,d,j,d}^{1D}(x_{t,d})$$

- ▶ Different number of addends (\rightarrow deferred to memory management)
- ▶ Small inner loop

Workload of one thread (index t):

Operation B :

$$\langle B_{t,\cdot}, v \rangle = \sum_{j=1}^m v_j \prod_d \phi_{l_{t,d}, i_{t,d}}^{1D}(x_{j,d})$$

Operation B^T :

$$\langle B_{\cdot,t}, v \rangle = \sum_{j=1}^g \alpha_j \prod_d \phi_{l_{j,d}, i_{j,d}}^{1D}(x_{t,d})$$

- ▶ Different number of addends (\rightarrow deferred to memory management)
- ▶ Small inner loop
- ▶ Thread cannot keep associated values in registers

Workload of one thread (index t):

Operation B :

$$\langle B_{t,\cdot}, v \rangle = \sum_{j=1}^m v_j \prod_d \phi_{l_{t,d}, i_{t,d}}^{1D}(x_{j,d})$$

Operation B^T :

$$\langle B_{\cdot,t}, v \rangle = \sum_{j=1}^g \alpha_j \prod_d \phi_{l_{j,d}, i_{j,d}}^{1D}(x_{t,d})$$

- ▶ Different number of addends (\rightarrow deferred to memory management)
- ▶ Small inner loop
- ▶ Thread cannot keep associated values in registers
- ▶ Bad access pattern to global memory

Implementation

Compile Time Constant Dimension

```
//...  
for  $d \leftarrow 0..D - 1$  do  
     $p \leftarrow p \cdot \phi_{l_{t,d}, i_{t,d}}^{1D}(x_j, d)$   
end for  
//...
```

template \langle int D \rangle

//...

for $d \leftarrow 0..D - 1$ **do**

$p \leftarrow p \cdot \phi_{l_{t,d}, i_{t,d}}^{1D}(x_{j,d})$

end for

//...

Compiler
 \Rightarrow

//...

$p \leftarrow p \cdot \phi_{l_{t,0}, i_{t,0}}^{1D}(x_{j,0})$

\vdots

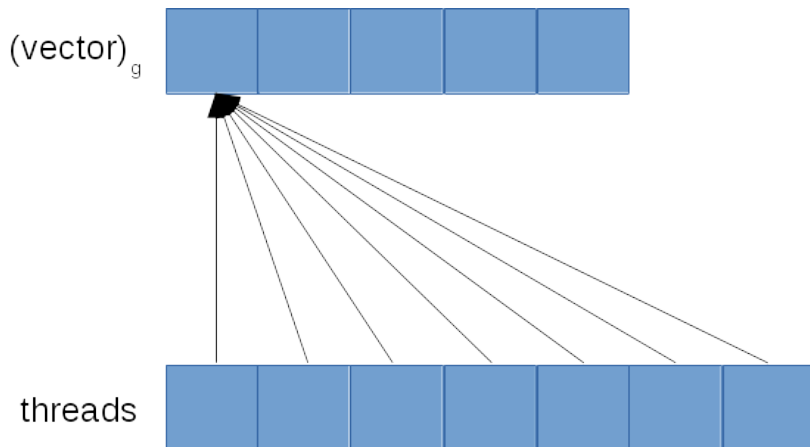
$p \leftarrow p \cdot \phi_{l_{t,D-1}, i_{t,D-1}}^{1D}(x_{j,D-1})$

//...

Implementation

Shared Memory

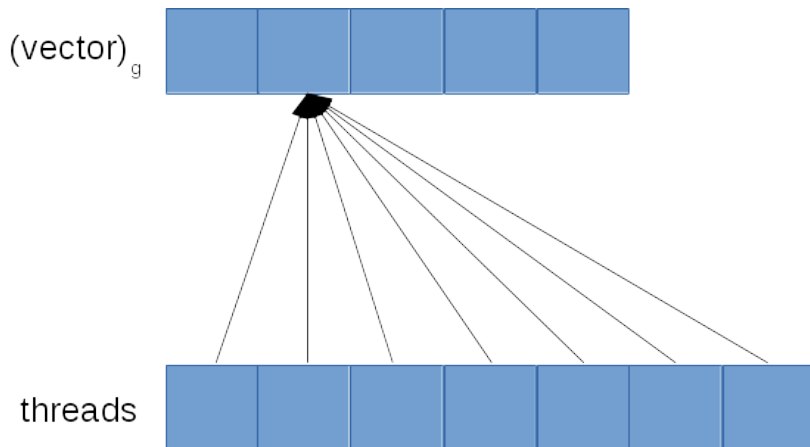
8 / 17



Implementation

Shared Memory

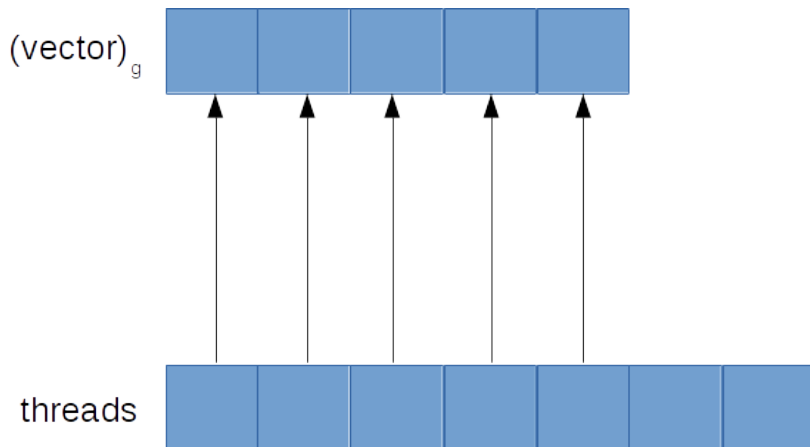
9 / 17



Implementation

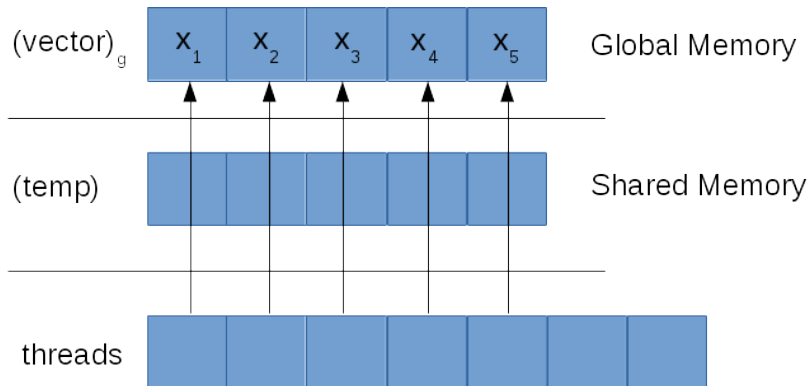
Shared Memory

10 / 17



Implementation

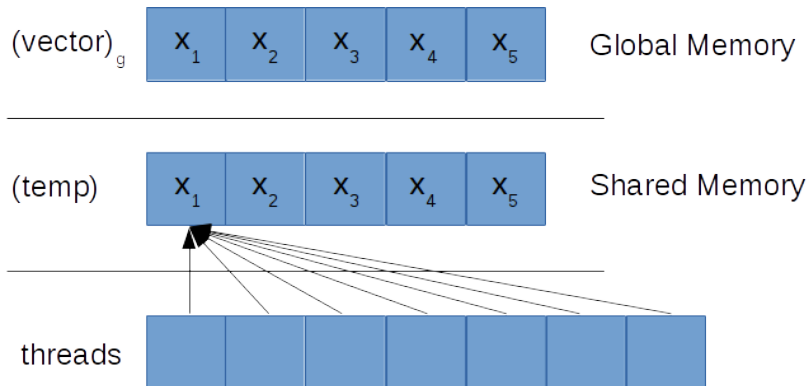
Shared Memory



Implementation

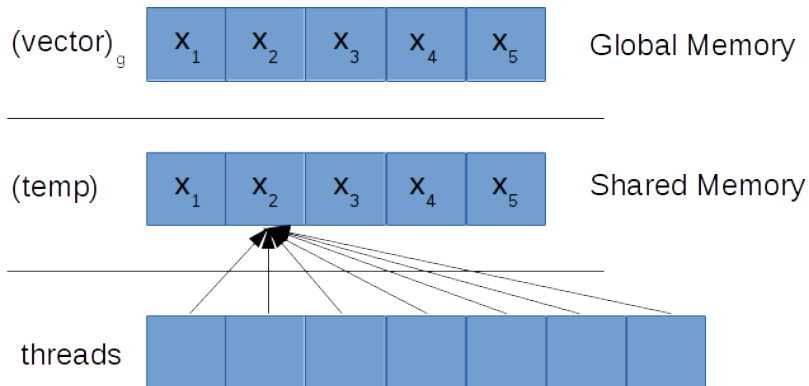
12 / 17

Shared Memory



Implementation

Shared Memory



The necessary data:

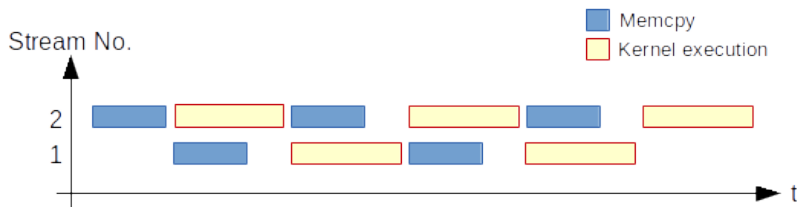
$$(l_{k,l}) \in \mathbb{R}^{g \times d}, (i_{k,l}) \in \mathbb{R}^{g \times d}, (x_{i,d}) \in \mathbb{R}^{m \times d}$$

Memory Management Problems:

- ▶ $m \gg g$ (Different number of addends)
- ▶ Program spends much time simply copying data
- ▶ Data might not fit entirely into device memory

Strategy:

1. Divide dataset into N chunks
2. Start $n < N$ streams and schedule the N chunks



- ▶ Single precision!
- ▶ DR5: Real-world dataset \sim 40 MB
- ▶ Intel Core i7 2600: **468.269 s**
- ▶ NVIDIA GeForce GTX680:

Templates	Streams	Shared Mem	Timing
X	X	X	812.242 s
✓	X	X	156.367 s
✓	✓	X	94.2451 s
✓	✓	✓	77.9071 s

- ▶ Discretization via Sparse Grids
- ▶ Least-Squares-Problem; minimization yields SLE
- ▶ Important techniques:
 - ▶ Allow the compiler to use registers
 - ▶ Unroll inner loops
 - ▶ Data access patterns + memory hierarchy
 - ▶ Balance operations
 - ▶ Streaming