# Parallelizing Regularized Least Squares

# Manuel A. Rivas

Massachusetts Institute of Technology

May 15, 2008

Manuel A. Rivas Parallelizing Regularized Least Squares

A (1) > A (1) > A

- Two weeks ago Prof. Edelman covered Support Vector Machines - classification type problems.
- Today, I will talk about my project on Regularized Least Squares - Ryan Rifkin a former student of Prof. Edelman is well known for his contribution to RLS.
- **③** Problem : Suppose we have pairs  $(x_i, y_i)$  with  $x_i \in \mathbb{R}^d$ , and we seek a d-dimensional regression function f(x).
- Solution : We setup the problem as

$$\min_{f\in\mathcal{H}}\left[\sum_{i=1}^{N}L(y_i,f(x_i))+\lambda J(f)\right],$$

where L(y, f(x)) is a loss function, J(f) is a penalty functional,  $\mathcal{H}$  is a space of functions on which J(f) is defined.

・ロト ・同ト ・ヨト ・ヨト

- Linear regression, linear discriminant analysis, logistic regression and separating hyperplanes all rely on a linear model.
- Extremely unlikely that the true function f(X) is actually linear in X.
- In regression problems,  $f(X) = \mathbb{E}(Y|X)$  will typically be nonlinear.
- Representing f(X) by a linear model is a convenient and at times appropriate approximation to avoid overfitting.

・ロト ・同ト ・ヨト ・ヨト

• Kernel functions commonly used:

- Define the Kernel matrix K to satisfy  $K_{ij} = k(X_i, X_j)$ .
- Given an arbitrary point X<sub>\*</sub>, k(X, X<sub>\*</sub>) is a column vector whose *i*th entry is k(X<sub>i</sub>, X<sub>\*</sub>).

▲ □ ▶ ▲ □ ▶ ▲

• In this project we use the setup :

$$\min_{f \in \mathcal{H}} \left[ \sum_{i=1}^{N} (y_i, f(x_i))^2 + \lambda J(f) \right]$$

- This minimizes the weighted sum of the *total* **square** loss. This loss function is great for regression type problems.
- Note: Support Vector Machines for classification use the following loss function

$$V = |1 - yf(x)|_+$$

or the Heavyside step function

$$V = \theta(-f(x)y).$$

We use Matlab Star-P to parallelize the data. Each processor is assigned a part of the kernel matrix to perform computation. Data Parallelism is exploited.

xtrain = linspace(-5,5,ntp\*p)'; %Uniformly spaced points
xtest = linspace(-5,5,nptest\*p)';%Uniformly spaced points
...
ps = xtrain\*diagonalmet\*xsup'; %xsup = xtrain; X diag(I) X
...
K = exp(-ps/2); % Training Kernel
...

c = (K+lambda\*eye(n))\y;

・ロト ・ 同ト ・ ヨト ・ ヨト - -

# Examples I





(a) n = 32 training points

(b) n = 2048 training points

<ロ> <同> <同> <同> < 同>

Figure: RLS solution for  $(\sin(\pi x))^3$ 

# Examples II



(日) (同) (三) (三)

#### Performance

Figure: Time needed to execute an iteration of RLS with training points = (32,64,128,256,512,1024,2048)



# Speedup

Figure: Speedup with number of processors = 1,2,4,6,8,12



Manuel A. Rivas Parallelizing Regularized Least Squares

Image: A (1) → A (

э

э

## Future Work

• As of now the limit on the number of datapoints my implementation handles is ' 5,000 datapoints.RLS does not need the QP steps.

# Future Work

- As of now the limit on the number of datapoints my implementation handles is ' 5,000 datapoints.RLS does not need the QP steps.
- SVM Classification is much more difficult to parallelize A constrained quadratic programming problem.

Image: A (1) → A (

# Future Work

- As of now the limit on the number of datapoints my implementation handles is ' 5,000 datapoints.RLS does not need the QP steps.
- SVM Classification is much more difficult to parallelize A constrained quadratic programming problem.
- Current literature implementation for SVMs use datasets of 500K to 1M data points requires ' 3 days with one machine.

Image: A (1) → A (

# Future Work

- As of now the limit on the number of datapoints my implementation handles is ' 5,000 datapoints.RLS does not need the QP steps.
- SVM Classification is much more difficult to parallelize A constrained quadratic programming problem.
- Current literature implementation for SVMs use datasets of 500K to 1M data points requires ' 3 days with one machine.
- Graf H.P., and Vapnik V. (2007) have recently proposed the Spread Kernel Support Vector Machine. Take advantage of the parallelization capability of decomposition algorithms. Idea is to parallelize Kernels and distribute training sets and perform QP independently.

・ロト ・同ト ・ヨト ・ヨト