Divide & Conquer Methods for Large-Scale Data Analysis

Inderjit S. Dhillon Dept of Computer Science UT Austin

International Conference on Machine Learning and Applications Detroit, MI Dec 4, 2014

Joint work with C.-J. Hsieh and S. Si

Inderjit S. Dhillon Dept of Computer Science UT Austin Divide & Conquer Methods for Big Data

Machine Learning Applications







gene-gene network

fMRI



Image classification

Spam classification

Gmail -



How to develop machine learning algorithms that scale to massive datasets?

Inderjit S. Dhillon Dept of Computer Science UT Austin Divide & Conquer Methods for Big Data

Divide & Conquer Method

- Divide original (large) problem into smaller instances
- Solve the smaller instances
- Obtain solution to original problem by combining these solutions



(人間) ト く ヨ ト く ヨ ト

-

Divide & Conquer Method

Strategy:



 D_3

D₃₃

D₃₄

How to merge?

æ

D₃₂

Better

performance

 D_{31}

<ロ> <同> <同> < 同> < 同>

Three Examples in Machine Learning

classification using kernel SVM dimensionality reduction for social network analysis learning graphical model structure

Divide & Conquer Kernel SVM

Inderjit S. Dhillon Dept of Computer Science UT Austin Divide & Conquer Methods for Big Data

< ≣ > <

Linear Support Vector Machines (SVM)

• SVM is a widely used classifier.

Given:

- Training data points x_1, \cdots, x_n .
- Each $x_i \in \mathbb{R}^d$ is a feature vector:
- Consider a simple case with two classes: $y_i \in \{+1, -1\}$.
- Goal: Find a hyperplane to separate these two classes: if $y_i = 1$, $w^T x_i \ge 1 - \xi_i$; $y_i = -1$, $w^T x_i \le -1 + \xi_i$.



- Given training data $x_1, \dots, x_n \in \mathbb{R}^d$ with labels $y_i \in \{+1, -1\}$.
- SVM primal problem:

$$\min_{w,\xi} \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i$$

s.t. $y_i(w^T x_i) \ge 1 - \xi_i, \ \xi_i \ge 0, \ i = 1, \dots, n,$

- Given training data $x_1, \dots, x_n \in \mathbb{R}^d$ with labels $y_i \in \{+1, -1\}$.
- SVM primal problem:

$$\begin{split} \min_{w,\xi} \; & \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \; & y_i(w^T x_i) \geq 1 - \xi_i, \; \xi_i \geq 0, \; i = 1, \dots, n, \end{split}$$

• SVM dual problem:

$$\min_{\alpha} \frac{1}{2} \alpha^{T} Q \alpha - e^{T} \alpha,$$

s.t. $0 \le \alpha_{i} \le C$, for $i = 1, ..., n_{i}$

where $Q_{ij} = y_i y_j x_i^T x_j$ and $e = [1, ..., 1]^T$.

• At optimum: $w = \sum_{i=1}^{n} \alpha_i y_i x_i$,

• What if the data is not linearly separable?

< ∃ > <

• What if the data is not linearly separable?

$$\begin{array}{cccc} & & & & & & & & \\ & & & & & & & \\ & & & & & \\ &$$

Solution: map data x_i to higher dimensional(maybe infinite) feature space $\varphi(x_i)$, where the classes are linearly separable.

- Kernel trick: $K(x_i, x_j) = \varphi(x_i)^T \varphi(x_j)$.
- Various types of kernels:
 - Gaussian kernel: $K(x, y) = e^{-\gamma ||x-y||_2^2}$;
 - Linear kernel: $K(x, y) = x^T y$;
 - Polynomial kernel: $K(x, y) = (\gamma x^T y + c)^d$.
- Solve the dual problem for SVM.

Linear vs. Nonlinear SVM

• UCI Forest Cover dataset (transform to binary version): Predicting forest cover type from cartographic features:

(Elevation, slope, soil type, ...)

464,810 training samples, 54 features

Transform to binary problem: Lodgepole Pine vs. the rest 6 types

	Training Time	Prediction Time	Prediction Accuracy	
Liblinear	Liblinear 3.6 secs		76.35%	
LibSVM	1 day	70060	96.15%	
(Gaussian kernel)	I Udy	10002X		



Kernel SVM

- Recently, (Fernández-Delgado et al., 2014) evaluated 179 classifiers from 17 families using 121 UCI data sets.
- Kernel SVM is the second best classifier (slightly outperformed by random forest)

Rank	Acc.	κ	Classifier
32.9	82.0	63.5	parRF_t (RF)
33.1	82.3	63.6	rf_t (RF)
36.8	81.8	62.2	svm_C (SVM)
38.0	81.2	60.1	svmPoly_t (SVM)
39.4	81.9	62.5	rforest_R (RF)
39.6	82.0	62.0	elm_kernel_m (NNET)
40.3	81.4	61.1	svmRadialCost_t (SVM)
42.5	81.0	60.0	svmRadial_t (SVM)
42.9	80.6	61.0	C5.0_t (BST)
44.1	79.4	60.5	avNNet_t (NNET)

Table from (Fernández-Delgado et al., 2014)

・ 同 ト ・ ヨ ト ・ ヨ ト …

Scalability for kernel SVM

- Challenge for solving kernel SVMs:
 - Space: *O*(*n*²);
 - Training Time: $O(n^3)$.
 - Prediction Time: $O(\bar{n}d)$
 - (\bar{n} : number of support vectors)
- Our solution: Divide-and-Conquer SVM (DC-SVM)

	Training Time	Prediction Time	Prediction Accuracy	
Liblinear	3.6 secs	1x	76.35%	
LibSVM	1 day	78862x	96.15%	
DC-SVM (early)	10 mins	308×	96.12%	
DC-Pred++	6 mins	18.8x	95.19%	

DC-SVM with a single level – divide step

- Partition α into k subsets $\{\mathcal{V}_1, \ldots, \mathcal{V}_k\}$.
- Solve each subproblem independently:

$$\begin{split} \min_{\boldsymbol{\alpha}_{(i)}} & \frac{1}{2} (\boldsymbol{\alpha}_{(i)})^T \boldsymbol{Q}_{(i,i)} \boldsymbol{\alpha}_{(i)} - \boldsymbol{e}^T \boldsymbol{\alpha}_{(i)}, \\ \text{s.t. } & 0 \leq \boldsymbol{\alpha}_{(i)} \leq \boldsymbol{C}, \end{split}$$

• Approximate solution for the whole problem:

$$\bar{\boldsymbol{\alpha}} = [\bar{\boldsymbol{\alpha}}_{(1)}, \dots, \bar{\boldsymbol{\alpha}}_{(k)}].$$



• Time complexity for k subproblems: $O(n^3) \rightarrow O(n^3/k^2).$



DC-SVM with a single level – conquer step

- Use $\bar{\alpha}$ to initialize a global coordinate descent solver.
- Converges quickly if

(1) $\|\bar{\alpha} - \alpha^*\|$ is small. (2) Support vectors in α^* are correctly identified in $\bar{\alpha}$.



• Question: how to find a partitioning that satisfies both?

Data Partitioning

Theorem 1

For a given partition π , the corresponding $ar{m{\alpha}}$ satisfies

$$0\leq f(ar{m lpha})-f(m lpha^*)\leq (1/2)C^2D(\pi),$$

where $D(\pi) = \sum_{i,j:\pi(x_i)\neq\pi(x_j)} |K(x_i, x_j)|$.

- Want a partition which
 - (1) Minimizes $D(\pi)$.
 - (2) Has balanced cluster sizes (for efficient training).
- Use kernel kmeans (but slow).
- Two step kernel kmeans:
 - Run kernel kmeans on a subset of samples with size $m \ll n$.

伺 ト イ ヨ ト イ ヨ

• Identify the clusters for the rest of data.

Demonstration of the bound

- Covertype dataset with 10000 samples and $\gamma = 32$ (best in cross validation).
- Our data partition scheme leads to a good approximation to the global solution α^{*}.



Quickly identifying support vectors

• Divide-and-conquer scheme helps to quickly identify support vectors.

Theorem 2

If $\bar{\alpha}_i = 0$ and

$$abla_iar{f}(ar{lpha}) > \textit{CD}(\pi)(1+\sqrt{n}\textit{K}_{max}/\sqrt{\sigma_n\textit{D}(\pi)}),$$

where $K_{max} = \max_i K(x_i, x_i)$, then x_i will not be a support vector of the whole problem (i.e., $\alpha_i^* = 0$).



Inderjit S. Dhillon Dept of Computer Science UT Austin

Divide & Conquer Methods for Big Data

- A tradeoff in choosing k (number of partitions):
 - Small $k \Rightarrow$ smaller $\|\bar{\alpha} \alpha^*\|$ but more time to solve subproblems.
 - Large $k \Rightarrow$ larger $\|\bar{\alpha} \alpha^*\|$ but less time to solve subproblems.

• A tradeoff in choosing k (number of partitions):

- Small $k \Rightarrow$ smaller $\|\bar{\alpha} \alpha^*\|$ but more time to solve subproblems.
- Large $k \Rightarrow$ larger $\|\bar{\alpha} \alpha^*\|$ but less time to solve subproblems.

• Therefore we run DC-SVM with multiple levels.



- A tradeoff in choosing k (number of partitions):
 - Small $k \Rightarrow$ smaller $\|\bar{\alpha} \alpha^*\|$ but need more time to solve subproblems.
 - Large $k \Rightarrow$ larger $\|ar{lpha} lpha^*\|$ but less time to solve subproblems.
- Therefore propose to run DC-SVM with multiple levels.



- A tradeoff in choosing k (number of partitions):
 - Small $k \Rightarrow$ smaller $\|\bar{\alpha} \alpha^*\|$ but need more time to solve subproblems.
 - Large $k \Rightarrow$ larger $\|ar{lpha} lpha^*\|$ but less time to solve subproblems.
- Therefore propose to run DC-SVM with multiple levels.



Solve the intermediate level problems.

• A tradeoff in choosing k (number of partitions):

- Small $k \Rightarrow$ smaller $\|\bar{\alpha} \alpha^*\|$ but need more time to solve subproblems.
- Large $k \Rightarrow$ larger $\|\bar{\alpha} \alpha^*\|$ but less time to solve subproblems.

• Therefore propose to run DC-SVM with multiple levels.

{1,...,n}

Solve the original problem.

Early Prediction

• Prediction using the I-th level solution

faster training time; the prediction accuracy is close to or even better than the global SVM solution.

- Prediction: sign $(\sum_{i \in \mathcal{V}_{\pi(x)}} y_i \alpha_i K(x_i, x))$.
- Prediction time reduced from O(d(#SV)) to O(d(#SV)/k)
- Covtype (*k* = 256): 78862x → 308x
- Still not fast enough for real-time applications!

Illustrative Toy Example

Two Circle Data: each circle is a class; separable by kernel kmeans.



Inderjit S. Dhillon Dept of Computer Science UT Austin Divide & Conquer Methods for Big Data

Illustrative Toy Examples

Two Circle Data: each circle is a class; separable by kernel kmeans.



Illustrative Toy Examples

Two Circle Data: not separable by kernel kmeans



Illustrative Toy Examples

Two Circle Data: not separable by kernel kmeans



Methods included in comparisons

- $\bullet~\mathrm{DC}\text{-}\mathrm{SVM}\text{:}$ our proposed method for solving the global SVM problem.
- DC-SVM (EARLY): our proposed method with early stopping (at 64 clusters).
- LIBSVM (Chang and Lin, 2011)
- CASCADE SVM (Graf et al., 2005)
- FASTFOOD (Le et al., 2013)
- LASVM (Bordes et al., 2005)
- LLSVM (Zhang et al., 2012)
- SPSVM (Keerthi et al., 2006)
- LTPU (Moody and Darken., 1989)
- BUDGETED SVM (Wang et al., 2012; Djuric et al., 2013)
- AESVM (Nandan et al., 2014)

Results with Gaussian kernel.

	web	spam	covtype		mnist8m	
	$n = 2.8 \times 10^5, d = 254$		$n = 4.65 \times 10^5, d = 54$		$n = 8 \times 10^6, d = 784$	
	$C = 8, \gamma = 32$		$C = 32, \gamma = 32$		$\mathcal{C}=1,\gamma=2^{-21}$	
	time(s)	acc(%)	time(s)	acc(%)	time(s)	acc(%)
DC-SVM (early)	670	99.13	672	96.12	10287	99.85
DC-SVM	10485	99.28	11414	96.15	71823	99.93
LIBSVM	29472	99.28	83631	96.15	298900	99.91
LaSVM	20342	99.25	102603	94.39	171400	98.95
CascadeSVM	3515	98.1	5600	89.51	64151	98.3
LLSVM	2853	97.74	4451	84.21	65121	97.64
FastFood	5563	96.47	8550	80.1	14917	96.5
SpSVM	6235	95.3	15113	83.37	121563	96.3
LTPU	4005	96.12	11532	83.25	105210	97.82
Budgeted SVM	2194	98.94	3839	87.83	29266	98.8
AESVM	3027	98.90	3821	87.03	16239	96.6

・ロト ・回ト ・ヨト ・ヨト

æ

Results with Gaussian kernel

Time vs Prediction Accuracy



A ►

Results with Gaussian kernel



covtype prediction accuracy



MNIST8m prediction accuracy

Image: A image: A

э

- ∢ ⊒ →

Results with grid of C, γ

• Total time on the grid of parameters C, γ :

•
$$C = 2^{-10}, 2^{-6}, 2^1, 2^6, 2^{10}$$

• $\gamma = 2^{-10}, 2^{-6}, 2^1, 2^6, 2^{10}$



э

Fast Prediction for Kernel SVM

Inderjit S. Dhillon Dept of Computer Science UT Austin Divide & Conquer Methods for Big Data

() <) <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <)
 () <

э

Prediction for SVM

- Linear SVM: $y = sign(w^T x + b); O(d)$.
- Kernel SVM: $y = sign(\sum_{i=1}^{n} \alpha_i K(x, x_i)); O(d(\#SV)).$
- DC-SVM with early prediction: $y = sign(\sum_{i \in V_{\pi(x)}} \alpha_i K(x, x_i));$ O(d(#SV/k)).
- DC-Pred++: O(dt); $t \ll \#SV/k$

Goal: achieve kernel SVM's accuracy using linear SVM prediction time.

DC-SVM with Nyström Approximation

• DC-SVM with early prediction



Drawback: still too many support vectors in each block

• Use Nyström approximation within each block



Nyström Kernel Approximation

• Goal: rank-*m* approximation \tilde{G} to kernel matrix *G* based on landmark points u_1, \ldots, u_m .

$$G \approx \tilde{G} = CWC^T.$$

• Perform prediction on a new point x given the model α : $\sum \alpha_i \tilde{K}(x, x_i) = \tilde{x}^T W C^T \alpha = \tilde{x}^T \beta$

where
$$\tilde{x} = [K(x, u_1), \dots, K(x, u_m)]^T$$
.

- The main computation is to form \tilde{x} .
- Time complexity for prediction: O(md).

Prediction



Inderjit S. Dhillon Dept of Computer Science UT Austin Divide & Conquer Methods for Big Data

Improvements over Traditional Nyström Approximation

- Traditional Nyström Approximation emphasizes kernel approximation error $\|G \tilde{G}\|_{F}$.
- We are interested in model approximation error $\| oldsymbol{lpha}^* oldsymbol{ar{lpha}} \|.$
- Use α -weighted kmeans centroids as the landmark points.
- Add pseudo landmark points for faster prediction.

DC-Pred++

• Training:



• Prediction:



æ

- 4 同 6 4 日 6 4 日 6

Dataset	Metric	DC-Pred++	LDKL	kmeans Nyström	AESVM	Liblinear
Letter	Prediction Time	12.8x	29x	140×	1542x	1x
n = 12,000	Accuracy	95.90%	95.78%	87.58%	80.97%	73.08%
CovType	Prediction Time	18.8x	35x	200x	3157x	1x
n = 522,910	Accuracy	95.19%	89.53%	73.63%	75.81%	76.35%
USPS	Prediction Time	14.4x	12.01x	200x	5787x	1x
n = 7291	Accuracy	95.56%	95.96%	92.53%	85.97%	83.65%
Webspam	Prediction Time	20.5x	23x	200x	4375x	1x
n = 280,000	Accuracy	98.4%	95.15%	95.01%	98.4%	93.10%

・ロン ・部 と ・ ヨ と ・ ヨ と …

æ

Social Network Analysis

Inderjit S. Dhillon Dept of Computer Science UT Austin Divide & Conquer Methods for Big Data

A B M A B M

Social Networks

- Nodes: individual actors (people or groups of people)
- Edges: relationships (social interactions) between nodes

Example: Karate Club Network:



Multi-Scale Link Prediction

• Combine predictions at each level to make final predictions



伺 ト イヨト イヨト

э

Experimental Results

Flickr dataset: 1.9M users & 42M links.

- Test set: sampled 5K users.
- Precision at top-100 & AUC results:



- ● ● ●

*using network-based features.

Learning Graphical Model Structure

Inderjit S. Dhillon Dept of Computer Science UT Austin Divide & Conquer Methods for Big Data

Sparse Inverse Covariance Estimation

- Given: *n* i.i.d. samples $\{y_1, \ldots, y_n\}$, $y_i \sim \mathcal{N}(\mu, \Sigma)$,
- Goal: Estimate the inverse covariance $\Theta = \Sigma^{-1}$.
- The resulting optimization problem:

$$\Theta = \arg\min_{X \succ 0} \left\{ -\log \det X + \operatorname{tr}(SX) + \lambda \|X\|_1 \right\} = \arg\min_{X \succ 0} f(X),$$

where $||X||_1 = \sum_{i,j=1}^n |X_{ij}|$ and $S = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\mu})(y_i - \hat{\mu})^T$.

- Regularization parameter $\lambda > 0$ controls the sparsity.
- Real world example: graphical model which reveals the relationships between Senators: (Figure from Banerjee et al, 2008)



Inderjit S. Dhillon Dept of Computer Science UT Austin Divid

Divide & Conquer Methods for Big Data

Divide-and-Conquer QUIC

• How can we solve high dimensional problems?

Time complexity is $O(p^3)$, so computationally expensive to solve it directly.

• Solution: Divide & Conquer!







 Θ from level-1 clusters $~~\Theta$ from level-2 clusters.

Inderjit S. Dhillon Dept of Computer Science UT Austin Divide & Conquer Methods for Big Data

Performance of Divide & Conquer QUIC



Conclusions

- Divide & Conquer approach for big data problems
 - Divide & Conquer Kernel SVM fast training and fast prediction
 - Social Network Analysis: Multi-Scale Link Prediction
 - Sparse Inverse Covariance Estimation: Divide & Conquer QUIC
- Interesting questions
 - What other data analysis problems can benefit from a divide & conquer approach? Regression? Matrix completion? Multi-label learning?
 - What are the corresponding algorithms?
 - Theoretical/statistical guarantees?
 - Can one exploit the multiplicity of local and global models for better prediction?
 - Parallelization? Needs dynamic load balancing, asynchronous parallelism.

< 回 > < 回 > < 回 >

References

Divide and Conquer Kernel SVM

[1] C.-J. Hsieh, S. Si and I. S. Dhillon A Divide-and-Conquer Solver for Kernel Support Vector Machines, ICML, 2014.

[2] S. Si, C.-J. Hsieh and I. S. Dhillon Memory Efficient Kernel Approximation, ICML, 2014.

[3] C.-J. Hsieh, S. Si and I. S. Dhillon *Fast Prediction for Large-Scale Kernel Machines*, NIPS, 2014.

Divide and Conquer Link-Prediction/Eigen-decomposition

[4] D. Shin, S. Si and I. S. Dhillon Multi-Scale Link Prediction, CIKM, 2012.

[5] S. Si, D. Shin, I. S. Dhillon and B. Parlett *Multi-Scale Spectral Decomposition of Massive Graphs*, NIPS, 2014.

Divide and Conquer Sparse Inverse Covariance Estimation

[6] C.-J. Hsieh, M. Sustik I. S. Dhillon and P. Ravikumar *Sparse Inverse Covariance Matrix Estimation using Quadratic Approximation*. NIPS, 2011.

[7] C.-J. Hsieh, I. S. Dhillon, P. Ravikumar and A. Banerjee *A Divide-and-Conquer Method for Sparse Inverse Covariance Estimation*, NIPS, 2012.

[8] C.-J. Hsieh, M. A. Sustik, I. S. Dhillon, P. Ravikumar, and R. A. Poldrack *BIG & QUIC: Sparse Inverse Covariance Estimation for a Million Variables*, NIPS, 2013.

[9] C.-J. Hsieh, I. S. Dhillon, P. Ravikumar and S. Becker, and P. A. Olsen *QUIC & DIRTY: A Quadratic Approximation Approach for Dirty Statistical Models*, NIPS, 2014.

Inderjit S. Dhillon Dept of Computer Science UT Austin Divide & Conquer Methods for Big Data