

Sublinear-Time Learning for High-Dimensional Models with Large Output Domains

Ian En-Hsu Yen

Machine Learning Department
Carnegie Mellon University

February 20, 2018

Outline

1 Inherent Sparsity in High-Dimensional Models

- Extreme Classification
- Sparsifying Neural Network

2 Decomposition for Sublinear-Time Learning

- Extreme Classification via Data Structures
- Structured Prediction of Large Output Domain

3 Sparse Estimator for Latent-Variable Models

- Binary Matrix Factorization

Outline

1 Inherent Sparsity in High-Dimensional Models

- Extreme Classification
- Sparsifying Neural Network

2 Decomposition for Sublinear-Time Learning

- Extreme Classification via Data Structures
- Structured Prediction of Large Output Domain

3 Sparse Estimator for Latent-Variable Models

- Binary Matrix Factorization

Extreme Classification

- Many problems involve **large output domain** where **prediction** and **loss evaluation** are very expensive.

	K=1 class	K $\approx 3 \times 10^5$ Classes
Train Time	25 sec.	> 96 days
Model Size	2.7MB	$\approx 870\text{GB}$

Table: Linear classifier on Wiki-LSHTC.



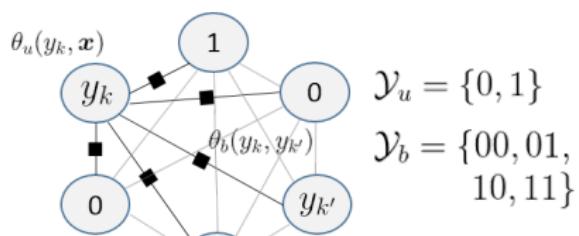
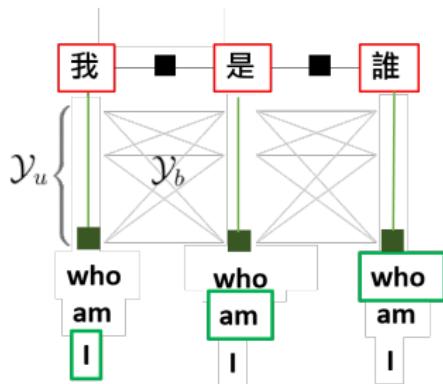
Wikipedia Dataset: $N \approx 10^6$, $K \approx 10^6$, $D \approx 10^6$



Megaface: $N \approx 10^6$, $K \approx 10^6$, $D \approx 10^3$

Extreme Structured Prediction

- Structured prediction and loss evaluation could have costs quadratic or cubic to the size of domain.



- Can we have training/prediction time, model size sublinear to the size of output domain K ?

Approaches to Large Output Domains

- **Structural Assumption:** low-rank, cluster hierarchy, binary encoding.
Issue: the structure might not holds \Rightarrow lower accuracy.

Approaches to Large Output Domains

- **Structural Assumption:** low-rank, cluster hierarchy, binary encoding.
Issue: the structure might not hold \Rightarrow lower accuracy.
- **Sampling Approximation:** sample a small portion of classes for the evaluation of loss and its gradient.
Issue: *classifier* confuses only few classes \Rightarrow slow convergence.

Approaches to Large Output Domains

- **Structural Assumption:** low-rank, cluster hierarchy, binary encoding.
Issue: the structure might not hold \Rightarrow lower accuracy.
- **Sampling Approximation:** sample a small portion of classes for the evaluation of loss and its gradient.
Issue: *classifier* confuses only few classes \Rightarrow slow convergence.
- **Question:** Can we speed it up *without approximation*?

Loss of Concentration (Dual Sparsity)

- $X : N \times D$ are features or embeddings of last layer in a NeuralNet.
- $W : D \times K$ are weights of K classes
- The empirical loss is $L_Y(\Theta) = \sum_{i=1}^N \ell_{y_i}(\Theta_{i,:})$, where $\Theta = XW$.
- Softmax Loss: $\ell_y(\theta) := \log(\sum_{k=1}^K \exp(\theta_k)) - \theta_y$
- Max-Margin Loss: $\ell_y(\theta) := [\max_{k_n \in \mathcal{N}_y} \theta_{k_n} - \min_{k_p \in \mathcal{P}_y} \theta_{k_p} + 1]_+$

$$-\nabla L_Y(\Theta^*) = \underbrace{\begin{matrix} k_A \text{ Support Labels} \\ \hline \text{---} & \text{---} & \text{---} \end{matrix}}_A \quad N \times K$$



Question: How to find those non-zero entries?

Approaches to Large Output Domains

- **Structural Assumption:** low-rank, clustered, hierarchy, binary encoding etc..
Issue: the structure might not hold \Rightarrow lower accuracy.
- **Sampling Approximation:** sample a small portion of classes for the evaluation of loss and its gradient.
Issue: classifier confuses only few classes \Rightarrow slow convergence.
- **Question:** Can we speed it up *without approximation*?
- **Learning by Search:** Find $k : \nabla_k L(\theta) \neq 0$.
 - ① Search via Sparse Matrix: $\theta = W^T x$.
 - ② Search via Data Structures: $\underset{k}{\operatorname{argmax}} \langle w_k, x \rangle$.

Joint Primal & Dual Sparsity

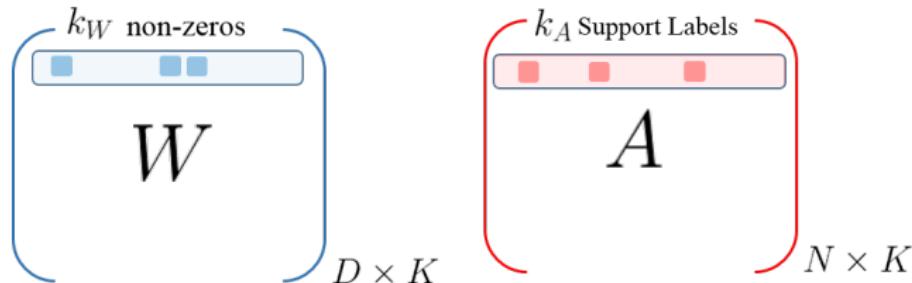
Theorem (Joint Primal & Dual Sparsity)

For any $\lambda > 0$ and X in general position,

$$W^* \in \underset{W}{\operatorname{argmin}} \lambda \|W\|_1 + L_Y(XW) \quad \text{and} \quad A^* \in -\nabla L_Y(\Theta^*)$$

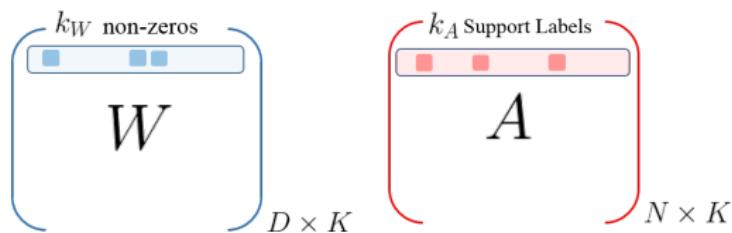
satisfy $Dk_W = \text{nnz}(W^*) \leq \text{nnz}(A^*) = Nk_A$.

- **Intuition:** Nk_A constraints vs Dk_W variables at optimal \Rightarrow under-determined.



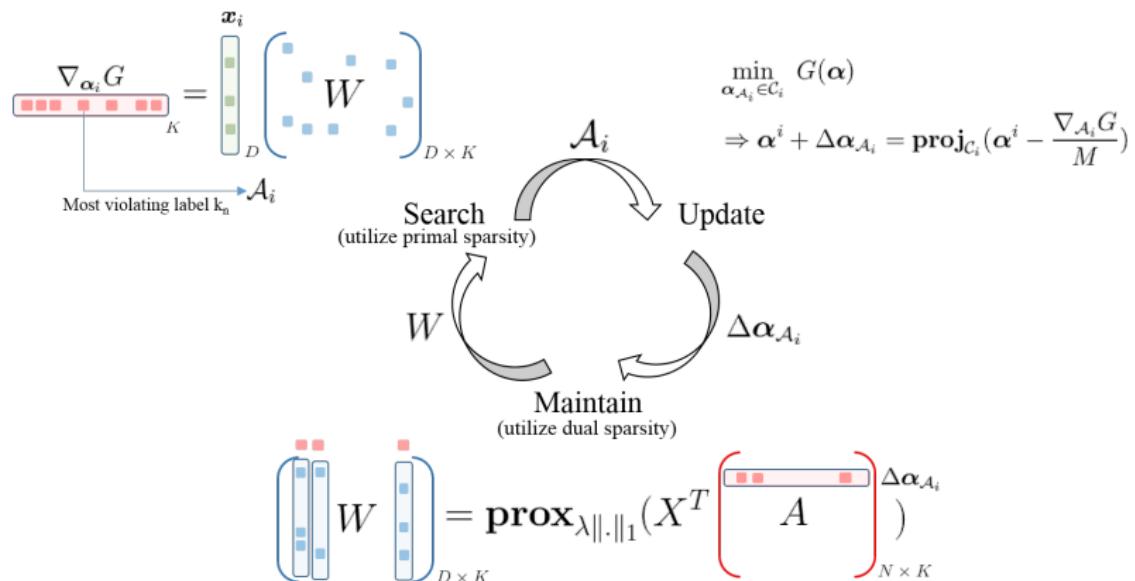
Search via Sparse Matrix

Data Set ¹	k_A : #support labels/sample	k_W : #nonzero/feature
bibtex (K=159)	18.17	1.94
aloi.bin (K=1,000)	3.24	0.31
EUR-Lex (K=3,956)	20.73	45.24
Dmoz (K=11,947)	5.87	0.116
LSHTC (K=12,294)	7.15	4.88
LSHTC-wiki (K=320,338)	18.24	20.95



¹Yen et al. "PD-Sparse: A primal and dual sparse approach to extreme multiclass and multilabel classification. ICML 16.

Dual Block-Coordinate Frank-Wolfe (BCFW)



- Sparse $W \Rightarrow$ Forward pass $\nabla L(W^T x)$ in $O(nnz(x)k_W)$.
- Sparse $\nabla L \Rightarrow$ Backward pass $\Delta W \sim \text{prox}_{\|\cdot\|_1}(-x \times \nabla L^T)$ in $O(nnz(x)k_A)$.
- $O(nnz(X)k_W + nnz(X)k_A)$ per epoch; $O(1/\epsilon)$ passes for ϵ -suboptimality.

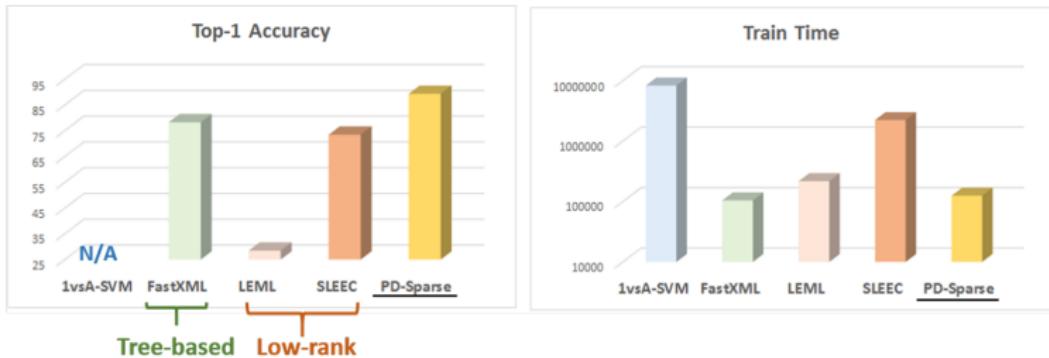
Results

LSHTC-wiki

$N = 2 \times 10^6$

$D = 2 \times 10^6$

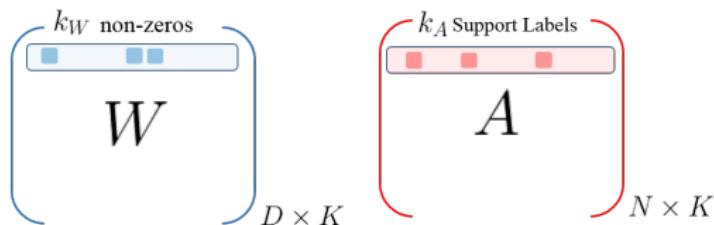
$K = 3 \times 10^5$



- **1-vs-All** takes > 3 months to train (single CPU).
- **PD-Sparse** takes ≈ 1 day to train. The accuracy is higher by $\approx 10\%$ than tree-based, low-rank, and local-embedding approaches.
- **PD-Sparse** predicts $\approx 100x$ faster than models from **1-vs-All**.
- (See more multiclass/multilabel results in (Yen et al. 2016).)

Algorithmic Choices

- **Dual BCFW**: better for maintaining **dual sparsity**.
- **Stochastic Gradient Descent**: better for **non-convex** objective.
- **Dual Greedy Coordinate Descent**: better for **parallelization**²
(solve each **column independently**; only applicable to **separable loss**.)



(Bounds on $\text{nnz}(\mathbf{W}^t)$, $\text{nnz}(\mathbf{A}^t)$ are algorithm-dependent.)

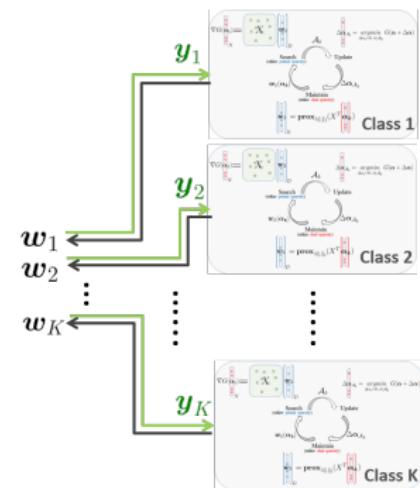
²Yen et al.. PPDSparse: A Parallel Primal-Dual Sparse Method for Extreme Classification. KDD 2017.

Algorithmic Choices

- **Dual BCFW**: better for maintaining **sparsity**.
- **Dual Greedy Coordinate Descent**: better for **parallelization** (solve each **column independently**; only applicable to **separable loss**).

Acc/Time	Wiki ($K=325k$)	Amazon ($K=670k$)
PPD-Sparse	64.1 / 6min	43.0 / 15min
PD-Sparse	60.7 / 1day	N/A
Parallel-1vsA	64.0 / 3day	43.0 / 2day
FastXML	57.2 / 6hr	33.1 / 1.5hr
SLEEC	58.3 / 11hr	35.6 / 6hr

Table: Parallel PD-Sparse (100 cores). ³



³Yen et al.. PPDSparse: A Parallel Primal-Dual Sparse Method for Extreme Classification. KDD 2017.

Outline

1 Inherent Sparsity in High-Dimensional Models

- Extreme Classification
- Sparsifying Neural Network

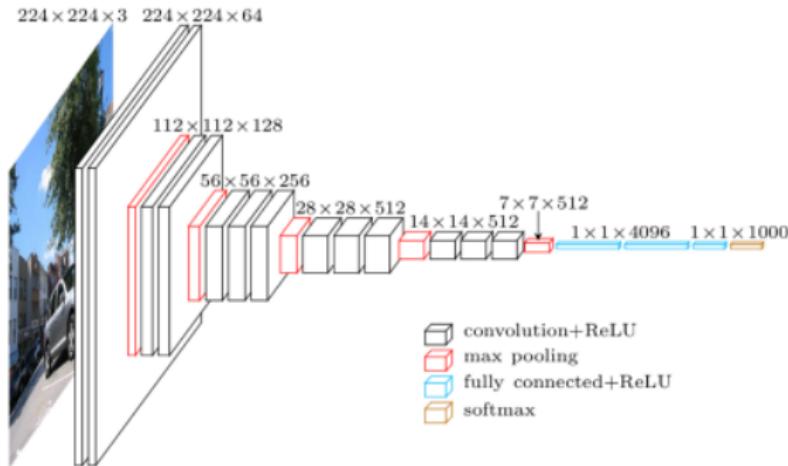
2 Decomposition for Sublinear-Time Learning

- Extreme Classification via Data Structures
- Structured Prediction of Large Output Domain

3 Sparse Estimator for Latent-Variable Models

- Binary Matrix Factorization

Sparsifying Neural Network



- A state-of-the-art DNN typically has $10^7 \sim 10^8$ parameters.
- **VGG16:** 138M parameters, with 102M parameters at the bottleneck layer.
- Many works try to sparsify DNNs, but how sparse? Is it data dependent?

Sparsifying Neural Network

Theorem (Sparse Neural Network)

Given a Neural Network $f_{W^*}(X)$ of loss L^* , for any $\epsilon > 0$ there is a sparse Neural Network $f_{\hat{W}}(X)$ of loss $L^* + \epsilon$ with

$$\text{nnz}(\hat{W}_{(l)}) \leq Nk_A, \quad l \in \text{layers}.$$

where k_A is the number of penalized classes per sample.

$\hat{W}_{(l)}$ can be any stationary point got by minimizing

$$\lambda \|W_{(l)}\|_1 + L(f_W(X)), \quad l \in \text{layers}$$

with initialization $W_{(l)}^*$ and some $\lambda(\epsilon) > 0$.

Sparsifying Neural Network

Theorem (Sparse Neural Network)

Given a Neural Network $f_{W^*}(X)$ of loss L^* , for any $\epsilon > 0$ there is a sparse Neural Network $f_{\hat{W}}(X)$ of loss $L^* + \epsilon$ with

$$\text{nnz}(\hat{W}_{(l)}) \leq Nk_A, \quad l \in \text{layers}.$$

where k_A is the number of penalized classes per sample.

Empirical Results: (without sacrificing accuracy)

- **ImageNet** ($N \approx 10^6$, $k_A \approx 5$):
 - $\dim(W) \approx 10^8 \Rightarrow \text{nnz}(\hat{W}) \approx 5 \times 10^6$ ($\approx 20x$ reduction).
- **CIFAR-10**: ($N = 5 \times 10^4$, $k_A \approx 0.1$):
 - $\dim(W) \approx 10^8 \Rightarrow \text{nnz}(\hat{W}) \approx 4 \times 10^4$ ($\approx 2000x$ reduction !!)
 - (current best known result: 68x)

Outline

1 Inherent Sparsity in High-Dimensional Models

- Extreme Classification
- Sparsifying Neural Network

2 Decomposition for Sublinear-Time Learning

- Extreme Classification via Data Structures
- Structured Prediction of Large Output Domain

3 Sparse Estimator for Latent-Variable Models

- Binary Matrix Factorization

Extreme Classification via Data Structures

- For DNNs of embedding size $10^2 \sim 10^3$, sparsity bound $\text{nnz}(\mathbf{W}^*) \leq Nk_A$ is not very useful ($DK \ll Nk_A$).
- Alternative:** Maximum Inner Product Search (MIPS):

$$k^*(\mathbf{x}) = \underset{k}{\operatorname{argmax}} [\mathbf{W}\mathbf{x}]_k = \underset{k}{\operatorname{argmax}} \mathbf{w}_k^T \mathbf{x}.$$

- Maintain vectors $\{\mathbf{w}_k\}_{k=1}^K$ in a data structure. Query by $\{\mathbf{x}_i\}_{i=1}^N$.
- Issue: high-recall MIPS is difficult in *high dimension*.

Extreme Classification via Data Structures

- **Alternative:** Maximum Inner Product Search (MIPS):

$$k^*(\mathbf{x}) = \underset{k}{\operatorname{argmax}} [W\mathbf{x}]_k = \underset{k}{\operatorname{argmax}} \mathbf{w}_k^T \mathbf{x}.$$

(Issue: high-recall MIPS is difficult in high dimension.)

- **To search accurately:** decompose high-dimensional MIPS into several low-dimensional MIPS with message passing.

$$\begin{pmatrix} \text{white} & \text{orange} & \text{white} & \text{orange} \\ W \end{pmatrix} \begin{pmatrix} \text{white} \\ \text{orange} \\ \text{white} \\ \text{orange} \\ \mathbf{x} \end{pmatrix} = W_{:,1}\mathbf{x}_1 + W_{:,2}\mathbf{x}_2 \dots + W_{:,F}\mathbf{x}_F$$

F factors: $\theta_1 + \theta_2 + \dots + \theta_F$

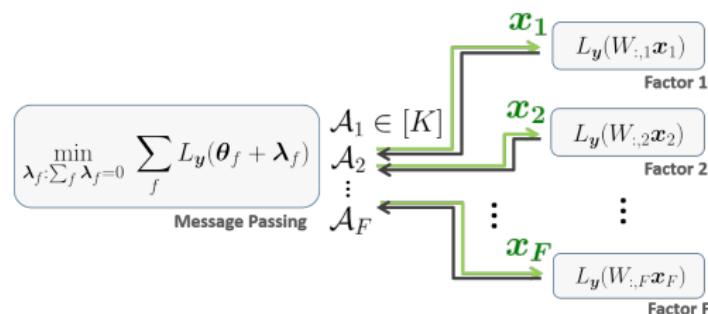
Extreme Classification via Data Structures

- **To search accurately:** decompose **high-dimensional MIPS** into several **low-dimensional MIPS** with message passing.

$$\begin{pmatrix} \text{blue} & \text{light orange} & \text{orange} \\ W \end{pmatrix} \begin{pmatrix} \text{light orange} \\ \text{orange} \\ \text{orange} \\ \text{orange} \\ \text{orange} \\ \text{orange} \end{pmatrix} = W_{:,1}\mathbf{x}_1 + W_{:,2}\mathbf{x}_2 \dots + W_{:,F}\mathbf{x}_F$$

F factors: $\theta_1 + \theta_2 + \dots + \theta_F$

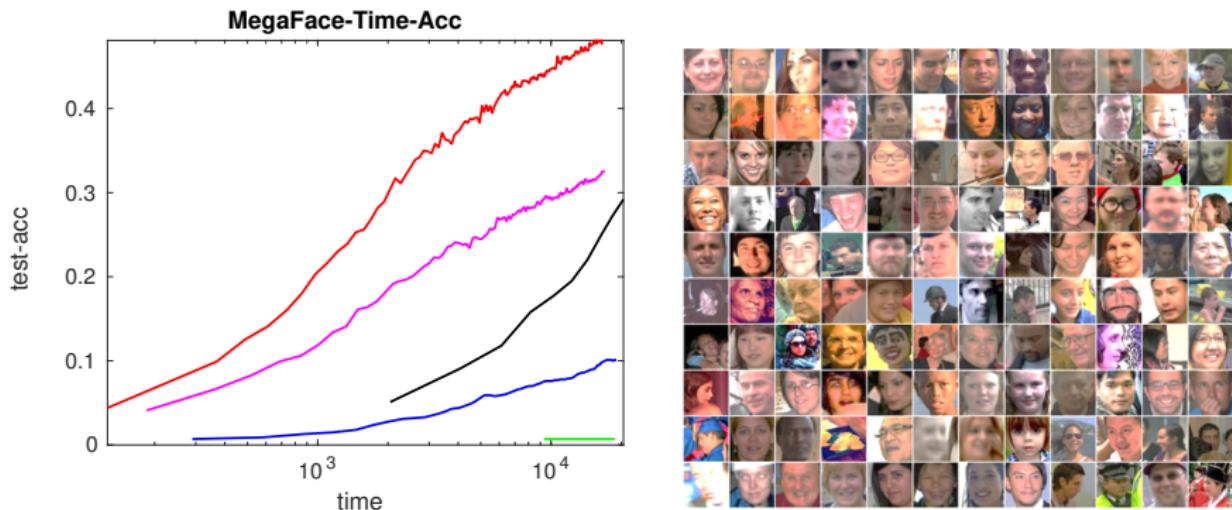
- **Key Insight:** $L(\sum_f \theta_f) = \min_{\lambda \in \Lambda} \sum_f L_f(\theta_f + \lambda_f)$.



(Messages λ_f are sparse due to sparse $\nabla L_f(\theta_f)$.)

Experiments: Face Recognition

Data set	#Person	#Image	Embedding Size	#Image / Person
MegaFace	672k	4.7M	128 (Facenet)	≈ 7



- Compare-1: **Max-Margin** vs. **MIPS** vs. **Dual-Decomp-MIPS**.
- Compare-2: **Max-Margin** vs. **Softmax** (and **Sampled-Softmax**).

Outline

1 Inherent Sparsity in High-Dimensional Models

- Extreme Classification
- Sparsifying Neural Network

2 Decomposition for Sublinear-Time Learning

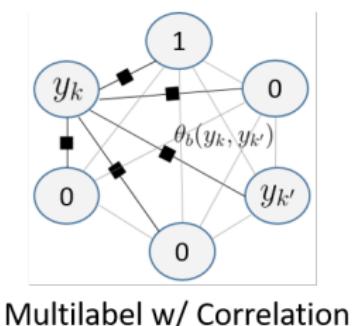
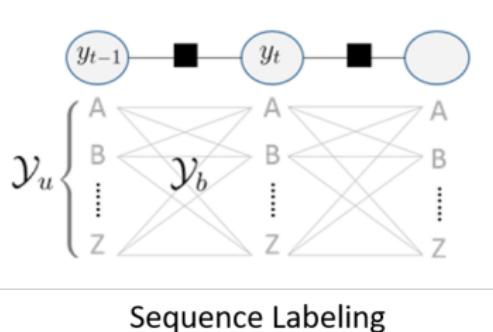
- Extreme Classification via Data Structures
- Structured Prediction of Large Output Domain

3 Sparse Estimator for Latent-Variable Models

- Binary Matrix Factorization

Structured Classifier with Large Output Domains

- What if there are multiple factors?
- Typically gradient evaluation takes $O(K^2)$ or more.

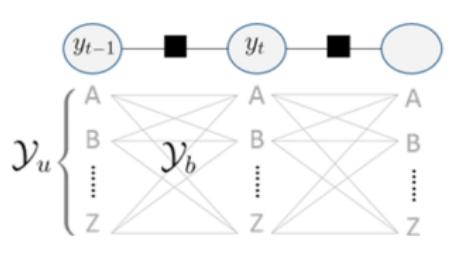


Structured Classifier with Large Output Domains

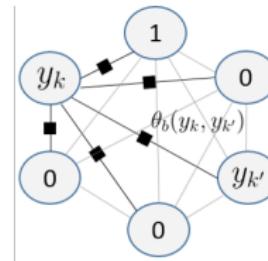
- Let $\theta(\mathbf{w}) := \sum_f \theta_f(\mathbf{w}) \in \mathbb{R}^{|\mathcal{Y}|}$ be the scores of each labeling $\mathbf{y} \in \mathcal{Y}$.
- Key Insight:**

$$\min_{\mathbf{w}} L \left(\sum_{f \in F} \theta_f(\mathbf{w}) \right) = \min_{\mathbf{w}} \min_{\lambda \in \Lambda} \sum_{f \in F} L_f \left(\theta_f(\mathbf{w}_f) + \sum_{j \in N(f)} \lambda_{jf} \right)$$

where λ_{jf} are messages from variable j to factor f .



Sequence Labeling



Multilabel w/ Correlation

Structured Classifier with Large Output Domains

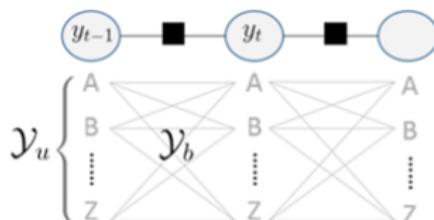
Key Insight: (λ_{jf} are messages from variable j to factor f .)

$$\underbrace{\min_{\mathbf{w}} \min_{\lambda \in \Lambda} \sum_{f \in F} L_f \left(\theta_f(\mathbf{w}_f) + \sum_{j \in N(f)} \lambda_{jf} \right)}_{\text{Inference}} = \underbrace{\min_{\lambda \in \Lambda} \min_{\mathbf{w}} \sum_{f \in F} L_f \left(\theta_f(\mathbf{w}_f) + \sum_{j \in N(f)} \lambda_{jf} \right)}_{\text{Multiclass Learning}}$$

- After solving **multiclass**, each factor's **#active classes** reduces from K to k_A , giving a sparse message λ_{jf} of size $O(k_A)$.

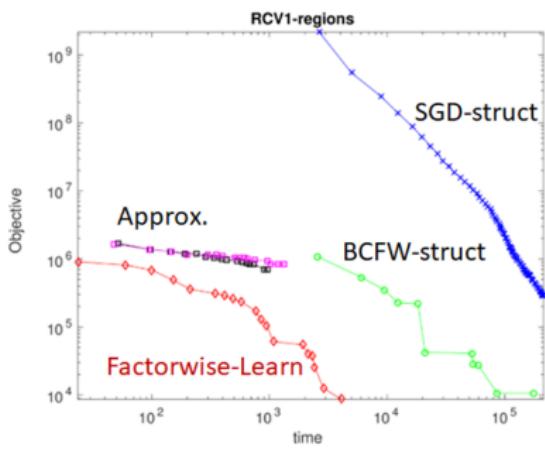
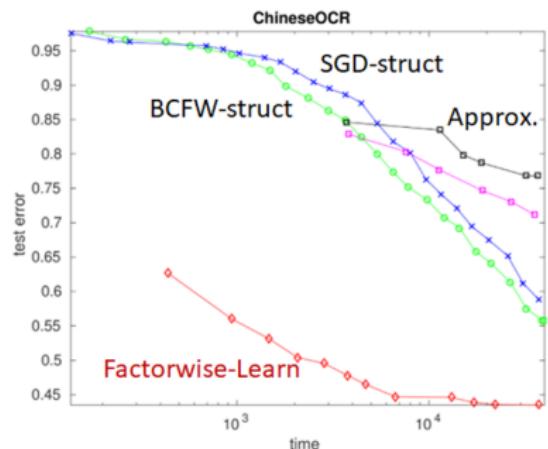
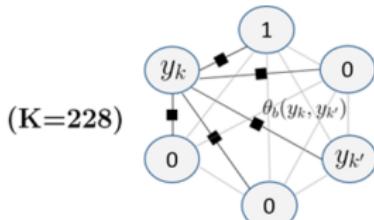
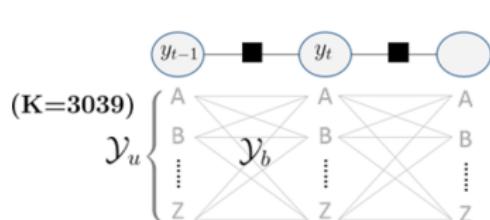
Inference: $O(NK^2 + NdK)$

multiclass: $O(Nk_A^2 + Nd \max(k_A, k_W))$



Sequence Labeling

Experiments: Sequence Labeling & Correlated Multilabel



³Yen. et al. Dual-Decomposed Learning with Factorwise Oracles for Structural SVMs of Large Output Domain. NIPS 2016.

Outline

1 Inherent Sparsity in High-Dimensional Models

- Extreme Classification
- Sparsifying Neural Network

2 Decomposition for Sublinear-Time Learning

- Extreme Classification via Data Structures
- Structured Prediction of Large Output Domain

3 Sparse Estimator for Latent-Variable Models

- Binary Matrix Factorization

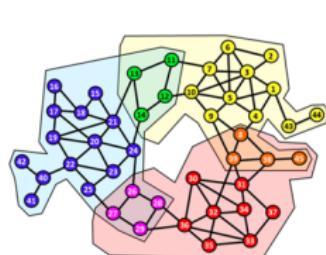
Binary Matrix Factorization (BMF)

- We are interested in finding factorization of the form:

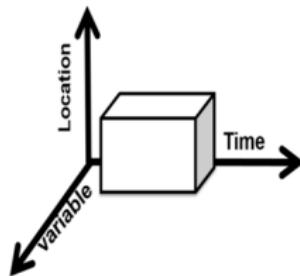
$$\min_{U,V} L(UV^T, R).$$

where $U \in \{0,1\}^{m \times k}$, and $V \in \{0,1\}^{n \times k}$ or $V \in \mathbb{R}^{n \times k}$.

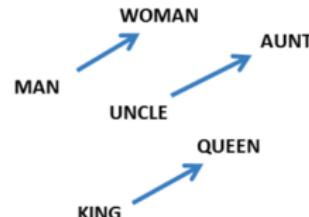
- Why binary? (interpretability, semi-supervision, and computational efficiency.)



Community Detection



Spatial-Temporal Factorization



Word Embedding

Binary Matrix Factorization (BMF)

- Let $U \in \{0, 1\}^{m \times k}$, and $V \in \{0, 1\}^{n \times k}$ or $V \in \mathbb{R}^{n \times k}$:

$$\min_{U, V} L(UV^T, R).$$

- Existing Approaches:**

- MCMC, Variational, Local Search (Indian Buffet Process):
No finite-time guarantee.
- Spectral Method:
Requires knowledge of **data distribution**.
- Algebraic Methods:
Can handle only exact case $R = UV^T$, in $O(mk2^k)$.

- Our Approach:**

- Relaxation in very high-dimensional space.
- First **polynomial-time** method of strong **approximation** guarantees.

BMF via Atomic Norm Minimization

- Let $\mathcal{S} := \{\mathbf{u}\mathbf{v}^T \mid \mathbf{u} \in \{0, 1\}^m, \mathbf{v} \in \{0, 1\}^n\}$.
- Define the [Atomic Norm](#):

$$\|M\|_{\mathcal{S}} := \min_{\mathbf{c} \geq 0} \sum_{\mathbf{u}\mathbf{v}^T \in \mathcal{S}} c_{\mathbf{u}, \mathbf{v}} \quad s.t. \quad M = \sum_{\mathbf{u}\mathbf{v}^T \in \mathcal{S}} c_{\mathbf{u}, \mathbf{v}} \mathbf{u}\mathbf{v}^T.$$

- Our [model estimator](#) is:

$$\min_M L_R(M) + \lambda \|M\|_{\mathcal{S}},$$

which is equivalent to:

$$\min_{\mathbf{c} \in \mathbb{R}_+^{|\mathcal{S}|}} L_R \left(\sum_{k=1}^{2^{m \times n}} c_k \mathbf{u}_k \mathbf{v}_k^T \right) + \lambda \|\mathbf{c}\|_1$$

Question: How to optimize with $|\mathcal{S}| = 2^{m \times n}$ variables?

Greedy Coordinate Descent via MAX-CUT

0. $\mathcal{A} = \emptyset$, $\mathbf{c} = 0$.

for $t = 1 \dots T$ **do**

1. Find an **approximate greedy** atom $\mathbf{u}\mathbf{v}^T$ via solving MAX-CUT:

$$\max_{\mathbf{u} \in \{0,1\}^m, \mathbf{v} \in \{0,1\}^n} \langle -\nabla L_R(M), \mathbf{u}\mathbf{v}^T \rangle.$$

2. Add $\mathbf{u}\mathbf{v}^T$ to an **active set** \mathcal{A} .

3. Refine $\mathbf{c}_{\mathcal{A}}$ via Proximal Gradient Method on:

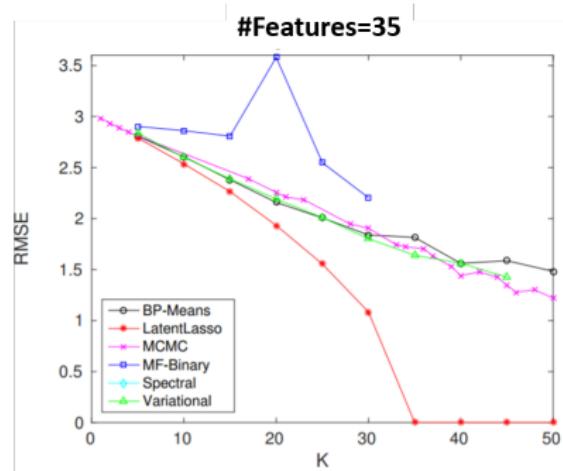
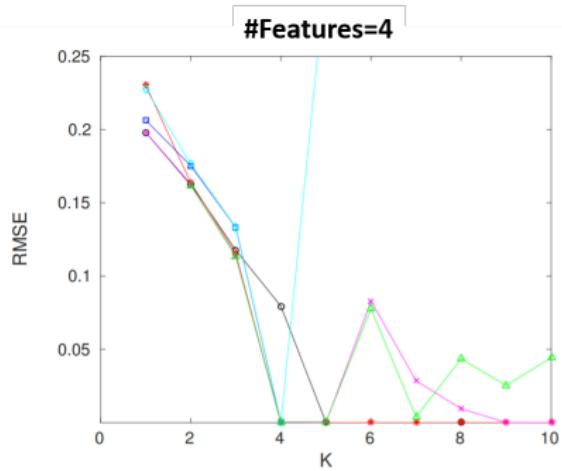
$$\min_{\mathbf{c} \geq 0} L_R\left(\sum_{k \in \mathcal{A}} c_k \mathbf{u}_k \mathbf{v}_k^T\right) + \lambda \|\mathbf{c}\|_1$$

4. Eliminate $\{\mathbf{u}_k \mathbf{v}_k^T | c_k = 0\}$ from \mathcal{A} .

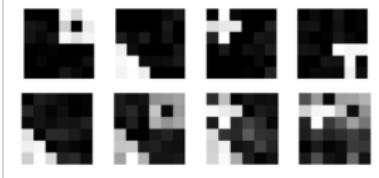
end for.

- MAX-CUT: 3/5-approximate solution in $O(nnz(\nabla L_R(M)))$.
- Running $\hat{k} = O(k^*/\epsilon)$ iterations, we have $L_R(\mathbf{U}\mathbf{S}\mathbf{V}^T) \leq L_R(\mathbf{U}^* \mathbf{S}^* \mathbf{V}^{*T}) + \epsilon$ for any binary factorization $\mathbf{U}^* \mathbf{S}^* \mathbf{V}^{*T}$ of rank k^* .

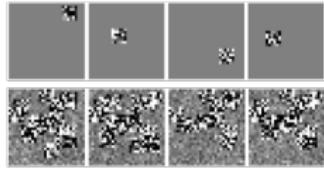
Results: Signal Decomposition Task



$K_{True}=4$



$K_{True}=35$



³See more results in: Yen et al. Latent Feature Lasso. ICML 2017.

Conclusion

- Many high-dimensional models have sparse structures, allowing us to:
 - Learn *Extreme Classifier* in sublinear time.
 - Sparsify Neural Network.
- Decomposition with sparse messages allows us to:
 - Handle large output domains for *Structured Predictor*.
 - Apply *Data Structures* effectively in learning.
- Ability to handle exponential number of variables allows us to design better estimators for some latent-variable models:
 - *Binary Matrix Factorization, Multiple Alignment, Motif Finding*, and *MAD-Bayes clustering*.

Appendix: Greedy Coordinate Descent via MAX-CUT

- The objective is:

$$\min_{\mathbf{c} \in \mathbb{R}_+^{|S|}} f(\mathbf{c}) + \lambda \|\mathbf{c}\|_1$$

- At each iteration, we find the coordinate of steepest descent:

$$j^* = \underset{j}{\operatorname{argmax}} -\nabla_j f(\mathbf{c}) = \underset{\mathbf{u} \in \{0,1\}^m, \mathbf{v} \in \{0,1\}^n}{\operatorname{argmax}} \langle -\nabla L_R(M), \mathbf{u} \mathbf{v}^T \rangle \quad (1)$$

that reduces to a problem similar to MAX-CUT:

$$\max_{\mathbf{u} \in \{0,1\}^m, \mathbf{v} \in \{0,1\}^n} \mathbf{u}^T C \mathbf{v} = \frac{1}{2} \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix}^T \begin{bmatrix} O & C \\ C & O \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix}$$

- It has **3/5-approximate** algorithm by rounding from a special type of SDP with $O(nnz(C))$ iterative solver.

Appendix: Binary Factor + Real-Valued Factor

- What if $V \in \mathbb{R}^{n \times k}$?

$$\min_{U \in \{0,1\}^{m \times k}} \left\{ \min_{V \in \mathbb{R}^{n \times k}} L_R(UV^T) + \frac{\tau}{2} \|V\|_F^2 \right\},$$

- Given U , the **dual problem w.r.t. V** is:

$$\min_{M=UU^T \in \{0,1\}^{m \times m}} \underbrace{\left\{ \max_{A \in \mathbb{R}^{m \times n}} -\frac{1}{2\tau} \text{tr}(AA^T M) - L_R^*(-A) \right\}}_{\tilde{L}_R(M)}.$$

- Define the Atomic Set: $\mathcal{S} := \{\mathbf{u}\mathbf{u}^T \mid \mathbf{u} \in \{0,1\}^m\}$. We have estimator:

$$\min_M \tilde{L}_R(M) + \lambda \|M\|_{\mathcal{S}}.$$