

# Improving the Generalization Performance of Multi-class SVM via Angular Regularization

Jianxin Li <sup>1</sup>, Haoyi Zhou <sup>1</sup>, Pengtao Xie <sup>2,3</sup>, Yingchun Zhang <sup>1</sup>

<sup>1</sup> Beihang University, <sup>2</sup> Carnegie Mellon University, <sup>3</sup> Petuum Inc

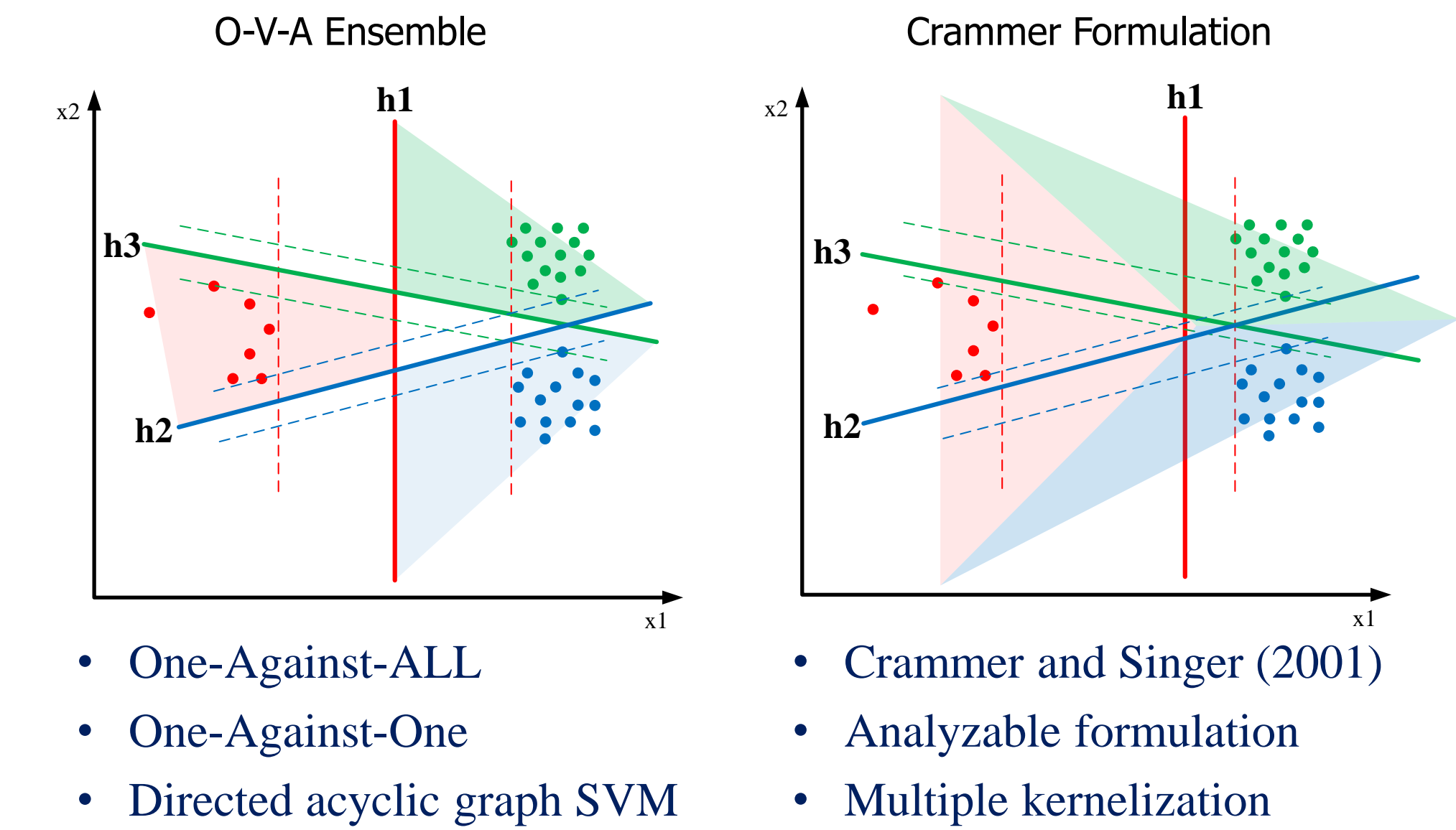
ICAI-17  
MELBOURNE

## I. ABSTRACT

In multi-class support vector machine (MSVM) for classification, one core issue is to regularize the coefficient vectors to reduce overfitting. Various regularizers have been proposed such as  $L_2$ ,  $L_1$ , and trace norm. In this paper, we introduce a new type of regularization approach -- angular regularization, that encourages the coefficient vectors to have larger angles such that class regions can be wider to flexibly accommodate unseen samples. We propose a novel angular regularizer based on the singular values of the coefficient matrix, where the uniformity of singular values reduces the correlation among different classes and drives the angles between coefficient vectors to increase. In generalization error analysis, we show that decreasing this regularizer effectively reduces generalization error bound. On various datasets, we demonstrate the efficacy of the regularizer in reducing overfitting.

## II. OVERVIEW

Multi-class SVM classification ( $K > 2$  classes):



Multi-class SVM formulation:

$$\min_{\mathbf{W}} P(\mathbf{W}) \triangleq \min_{\mathbf{W}} \frac{1}{m} \sum_{i=1}^m L(\mathbf{W}; (\mathbf{x}_i, y_i))$$

$$L(\mathbf{W}; (\mathbf{x}_i, y_i)) = \max(0, 1 + \mathbf{w}_{R_i}^T \cdot \mathbf{x}_i - \mathbf{w}_{y_i}^T \cdot \mathbf{x}_i),$$

where  $R_i = \operatorname{argmax}_{k \in \mathcal{Y}, k \neq y_i} \mathbf{w}_k^T \cdot \mathbf{x}_i$ .

Multi-class SVM with norm-based Regularization:

- L1-norm (Bradley et al, 1998; Wang et al, 2012),
- L2-norm (Weston et al, 1999; Guermuer, 2002),
- Doubly-norm (Wang et al, 2006)

$$\min_{\mathbf{W}} \lambda_1 \|\mathbf{W}\|_1 + \frac{1}{m} \sum_{i=1}^m L(\mathbf{W}; (\mathbf{x}_i, y_i))$$

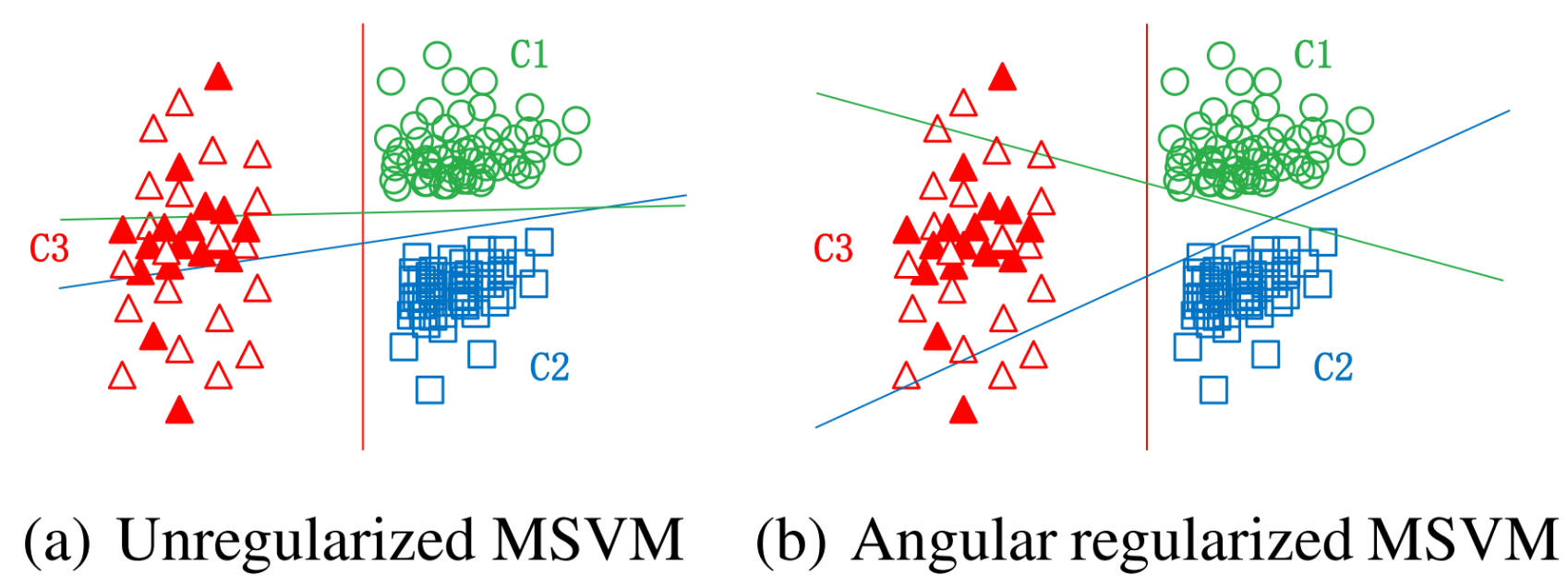
$$\min_{\mathbf{W}} \frac{\lambda}{2} \|\mathbf{W}\|_2^2 + \frac{1}{m} \sum_{i=1}^m L(\mathbf{W}; (\mathbf{x}_i, y_i))$$

$$\min_{\mathbf{W}} \lambda_1 \|\mathbf{W}\|_1 + \frac{\lambda_2}{2} \|\mathbf{W}\|_2^2 + \frac{1}{m} \sum_{i=1}^m L(\mathbf{W}; (\mathbf{x}_i, y_i))$$

Angular perspective to reduce Overfitting:

Table 1: Classification Accuracy of MSVM- $\ell_2$

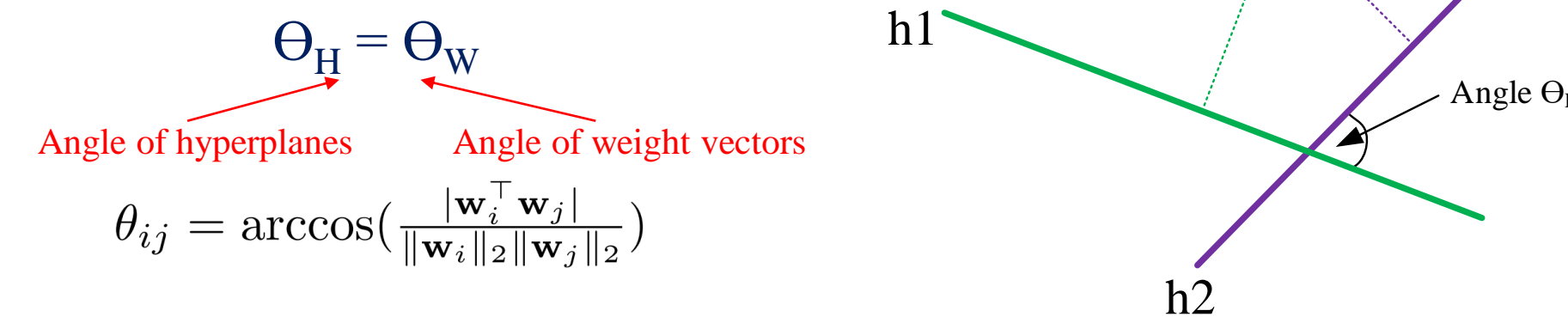
Dataset	Yeast	Usps	YaleB
Dim. of Features	8	256	1024
Train Accuracy (%)	57.23	97.31	97.70
Test Accuracy (%)	52.00	90.48	93.25



## III. METHODS

Desirable diversity measure:

INVARIANT to scale, translation, rotation and orientation of the two hyperplanes.



Angular Regularizer (AR):

The negative of the minimum angle:

$$\mathcal{R}(\mathbf{W}) = -\min_{i \neq j} \theta_{ij}$$

- Pros: directly encourages these coefficient vectors to have larger angles, and this definition facilitates theoretical analysis.
- Cons: non-smooth (similar to the one in Xie et al, 2015).

Decorrelation Regularizer (DR):

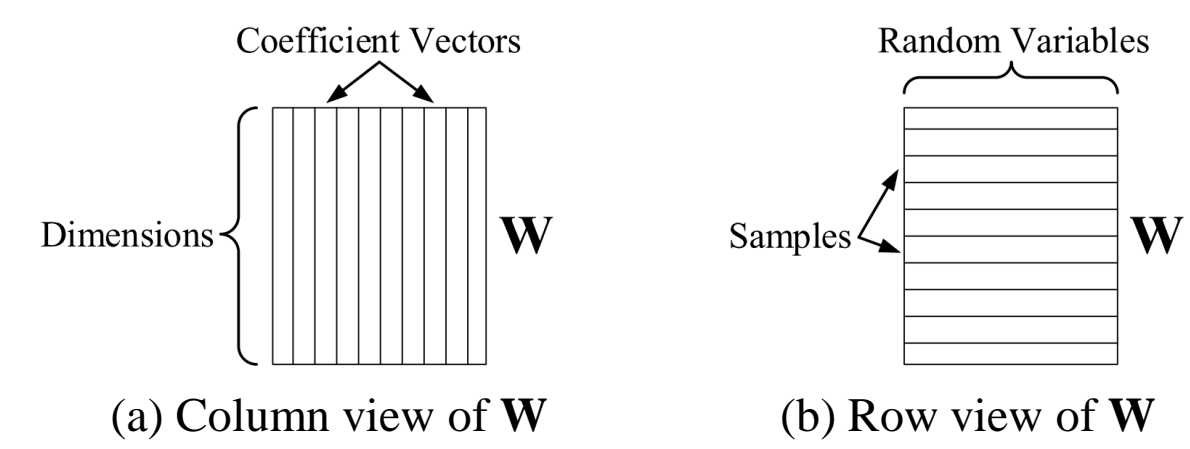
Given two hyperplanes ( $h_i$  and  $h_j$ ), we can treat  $i$  and  $j$  as two random variables:

$$\rho_{ij} = \frac{\sum_{n=1}^D (w_{in} - \bar{w}_i)(w_{jn} - \bar{w}_j)}{\sqrt{\sum_{n=1}^D (w_{in} - \bar{w}_i)^2 \sum_{n=1}^D (w_{jn} - \bar{w}_j)^2}}$$

$$\bar{w}_i = \frac{1}{n} \sum_{n=1}^D w_{in} = 0 \quad \bar{w}_j = \frac{1}{n} \sum_{n=1}^D w_{jn} = 0$$

$$\rho_{ij} = \cos(\theta_{ij})$$

A set of hyperplanes' weights vectors  $\mathbf{W}$  can be reformulated as full-rank matrix with rank  $K < D$ :  $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$ .



We compute the full-rank  $K \times K$  Gram matrix  $\mathbf{G} = \mathbf{W}^T \mathbf{W}$ , and

$$\mathbf{G} = \sum_{i=1}^K \lambda_i \mathbf{u}_i \mathbf{u}_i^T$$

Eigenvector  $\mathbf{u}_i$  of  $\mathbf{G}$  represents a principal direction of the point cloud.

The associated eigenvalue  $\lambda_i$  measures the variability of points along that direction.

(a)  $\lambda_1 > \lambda_2$  (b)  $\lambda_1 \approx \lambda_2$

We define a probability distribution of the eigenvalues

$$p(X = i) = \frac{\lambda_i}{\sum_{j=1}^K \lambda_j} \quad (\text{normalization})$$

If  $p(X)$  is close to uniform distribution  $q(X = i) = \frac{1}{K}$ , hyperplanes diversified. And we measure it by Kullback-Leibler divergence:

$$KL(q||p) = \sum_{i=1}^K \frac{1}{K} \log \frac{1/K}{\lambda_i / \sum_{j=1}^K \lambda_j}$$

$$= \log \sum_{j=1}^K \lambda_j - \frac{1}{K} \sum_{i=1}^K \log \lambda_i - \log K$$

Using two facts  $\operatorname{tr}(\mathbf{G}) = \sum_{j=1}^K \lambda_j$  and  $\det(\mathbf{G}) = \prod_{i=1}^K \lambda_i$ , we have

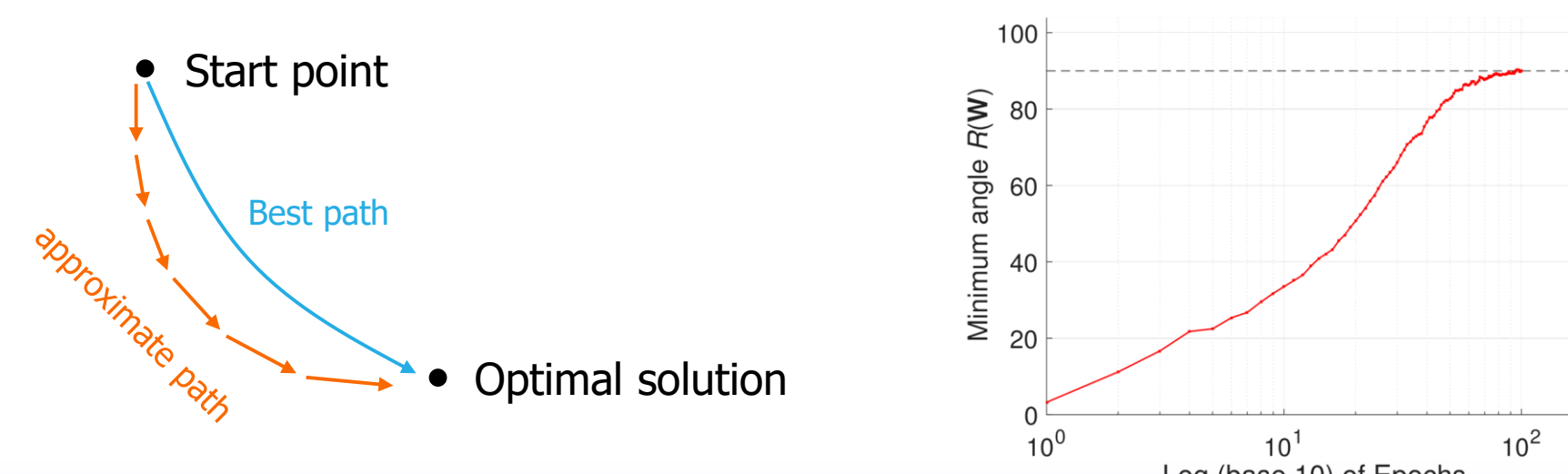
$$KL(q||p) = \log \operatorname{tr}(\mathbf{W}^T \mathbf{W}) - \frac{1}{K} \log \det(\mathbf{W}^T \mathbf{W}) - \log K$$

Dropped

The Decorrelation Regularizer is

$$\hat{\mathcal{R}}(\mathbf{W}) = \log \operatorname{tr}(\mathbf{W}^T \mathbf{W}) - \frac{1}{K} \log \det(\mathbf{W}^T \mathbf{W})$$

**Theorem 1.** Let  $\nabla \hat{\mathcal{R}}$  be the gradient of  $\hat{\mathcal{R}}(\mathbf{W})$  w.r.t  $\mathbf{W}$ .  $\exists \kappa > 0$ , such that  $\forall \eta \in (0, \kappa)$ ,  $\mathcal{R}(\mathbf{W} - \eta \nabla \hat{\mathcal{R}}) \leq \mathcal{R}(\mathbf{W})$ .



## IV. OPTIMIZATION

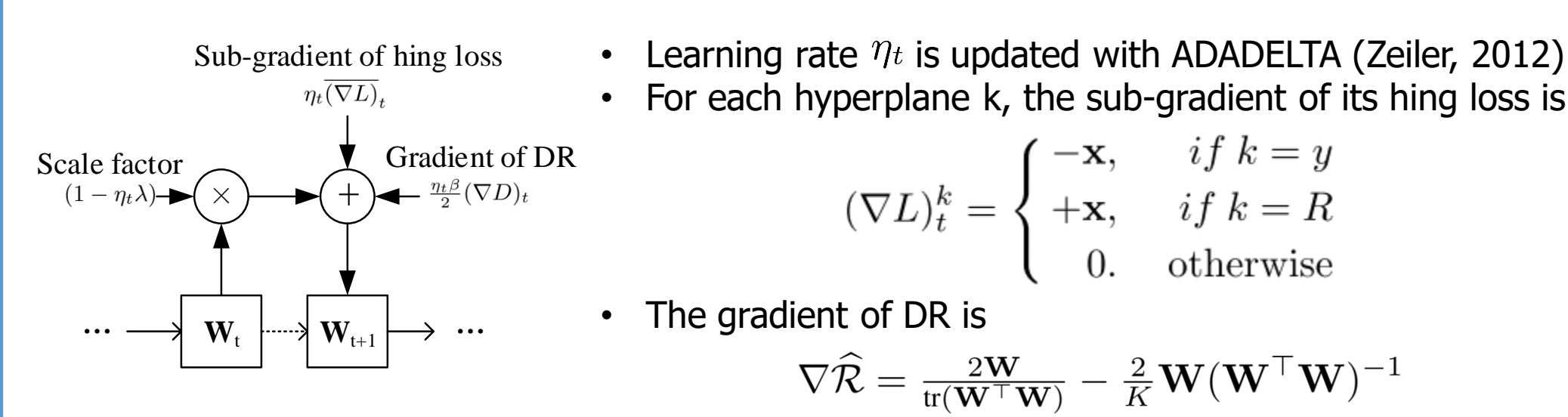
An Angle-regularized MSVM (AR-MSVM)

$$\min_{\mathbf{W}} \frac{1}{m} \sum_{i=1}^m L(\mathbf{W}; (\mathbf{x}_i, y_i)) + \frac{\lambda}{2} \|\mathbf{W}\|_2^2 + \beta \mathcal{R}(\mathbf{W})$$

Objective function for Optimization

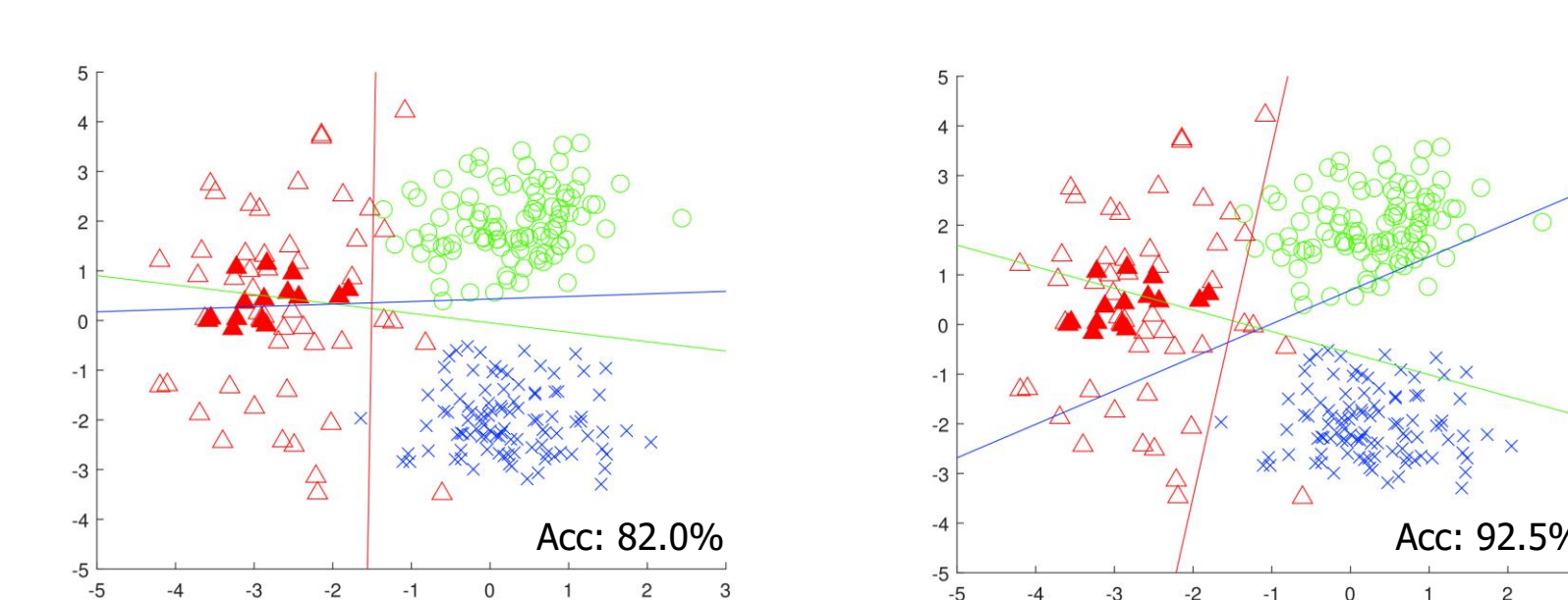
$$\min_{\mathbf{W}} \frac{1}{m} \sum_{i=1}^m \max(0, 1 + \mathbf{w}_{r_i}^T \mathbf{x}_i - \mathbf{w}_{y_i}^T \mathbf{x}_i) + \frac{\lambda}{2} \|\mathbf{W}\|_2^2 + \frac{\beta}{2} \hat{\mathcal{R}}(\mathbf{W})$$

Optimizing method: Pegasos-style Stochastic sub-gradient method (Shalev-Shwartz et al., 2011; Wang et al., 2010).



## V. EXPERIMENTAL RESULTS

Case study: Intuitive experiments



Multi-class Classification

Datasets:

Table 2: Statistics of Datasets

Dataset	#Classes	#Train	#Test	#Features
YaleB	38	1500	914	1024
ImageNet-50	50	50K	10K	128
Covtype	7	106K	40K	54
Shuttle	7	30450	14500	9
New-thyroid	3	108	107	5
Yeast	10	1134	350	8
Dermatology	6	323	35	33
Page-Blocks	5	4924	548	10
Wine-Quality-Red	6	1439	160	11
Zoo	7	89	12	16

Baseline:

- C-SVC (Chang and Lin, 2011): performing multi-class classification using an O-A-O strategy.
- AMM (Wang et al, 2011): a multi-class classification method based on adaptive multi-hyperplanes.
- MSVM-L1 (Shalev-Shwartz et al, 2011): L1-regularized MSVM.
- CS-MSVM (Yu et al, 2011): MSVM regularized by a cosine similarity regularizer.
- IC-MSVM (Bao et al, 2013): MSVM regularized by an incoherence regularizer.
- Div-MSVM (Xie et al, 2015): MSVM regularized by a diversity-promoting regularizer.
- MSVM-Struct (Tsochantaridis et al, 2004): L2-regularized MSVM with a cutting-plane algorithm.
- MSVM-L2 (Wang et al, 2010): L2-regularized MSVM solved with a SGD algorithm.

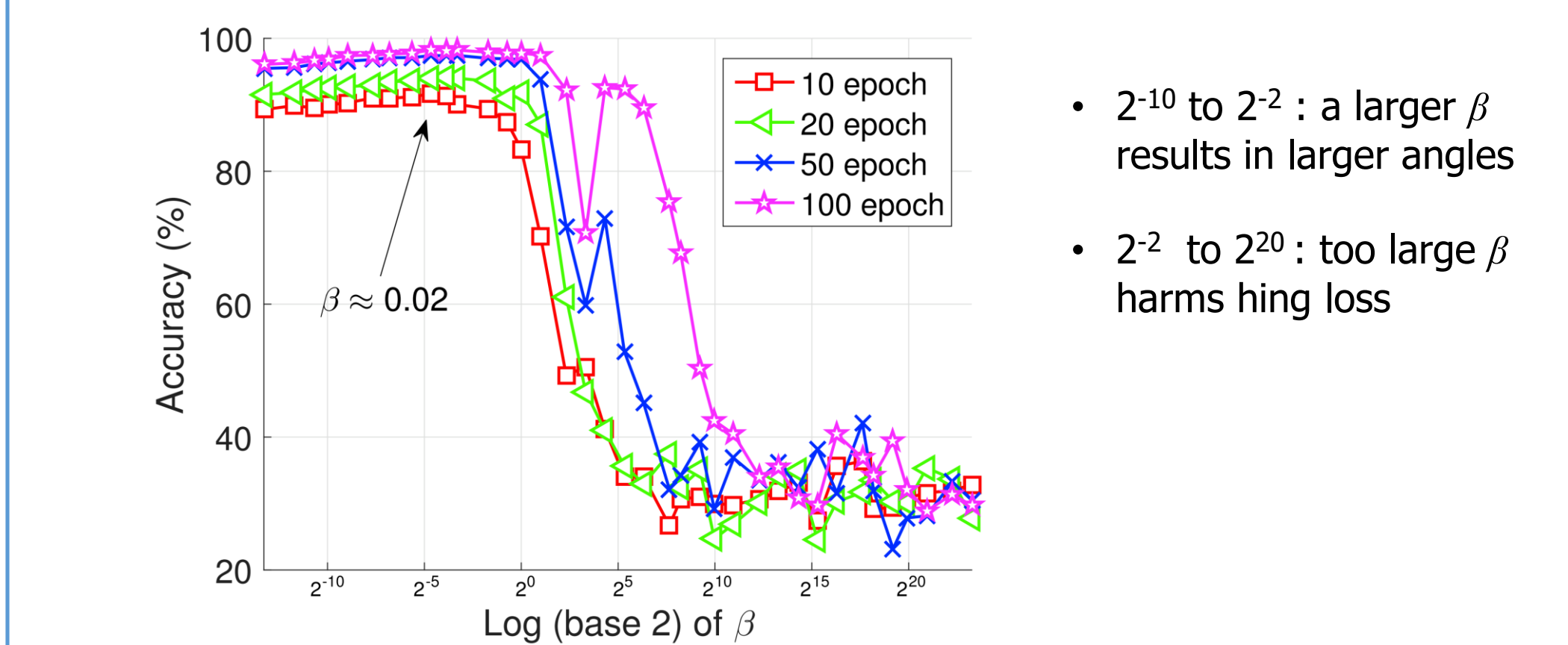
Table 3: Classification results (%) on six multi-class datasets

Dataset	Metric	C-SVC	AMM	MSVM-L1	CS-MSVM	IC-MSVM	Div-MSVM	MSVM-Struct	MSVM- $\ell_2$	AR-MSVM
YaleB	Acc	55.31 ± 0.1	83.32 ± 0.5	85.12 ± 0.4	91.13 ± 0.6	90.54 ± 0.9	94.25 ± 0.4	91.50 ± 0.0	93.90 ± 0.5	94.55 ± 0.8
	F	54.74 ± 0.0	86.88 ± 0.8	87.11 ± 0.6	93.71 ± 0.3	90.54 ± 0.4	93.87 ± 0.3	91.39 ± 0.0	93.51 ± 0.5	94.38 ± 0.7
ImageNet-50	Acc	91.59 ± 0.4	92.27 ± 0.1	91.23 ± 0.0	92.86 ± 0.1	92.96 ± 0.1	92.93 ± 0.1	89.36 ± 0.0	92.32 ± 0.1	93.12 ± 0.1
	F	91.70 ± 0.4	92.31 ± 0.1	91.38 ± 0.0	92.90 ± 0.1	92.98 ± 0.1	92.97 ± 0.1	89.65 ± 0.0	92.35 ± 0.0	93.17 ± 0.1
Covtype	Acc	71.35 ± 0.5	70.21 ± 0.3	56.12 ± 0.4	70.29 ± 0.1	70.82 ± 0.2	70.70 ± 0.1	65.84 ± 0.1	69.45 ± 0.2	71.75 ± 0.1
	F	50.97 ± 0.1	49.84 ± 0.4	37.03 ± 0.5	50.67 ± 0.1	50.82 ± 0.1	50.53 ± 0.3	36.80 ± 0.1	45.64 ± 0.8	51.31 ± 0.9
Shuttle	Acc	98.86 ± 0.3	93.31 ± 0.3	79.41 ± 0.4	96.02 ± 0.0	96.22 ± 0.1	94.74 ± 0.2	66.57 ± 0.0	92.83 ± 0.1	97.08 ± 0.3
	F	52.94 ± 0.1	54.12 ± 0.3	24.90 ± 0.0	57.20 ± 0.1	57.57 ± 0.1	55.61 ± 0.2	33.66 ± 0.0	53.55 ± 0.2	58.53 ± 0.3
New-thyroid	Acc	76.92 ± 0.9	86.77 ± 0.8	86.15 ± 0.0	78.46 ± 0.9	78.69 ± 0.3	91.54 ± 0.8	89.23 ± 0.0	90.87 ± 0.0	92.15 ± 0.5
	F	58.73 ± 0.0	83.66 ± 0.4	87.50 ± 0.0	75.32 ± 0.4	78.46 ± 0.1	90.08 ± 0.1	87.50 ± 0.0	89.24 ± 0.0	90.91 ± 0.7
Yeast	Acc	59.71 ± 0.5	48.63 ± 0.3	53.71 ± 0.4	53.61 ± 0.2	54.14 ± 0.2	53.14 ± 0.5	49.93 ± 0.1	51.86 ± 0.3	54.80 ± 0.1
	F	61.81 ± 0.0	44.87 ± 0.1	48.63 ± 0.4	52.80 ± 0.1	52.07 ± 0.0	52.64 ± 0.3	37.09 ± 0.1	49.45 ± 0.1	54.65 ± 0.4

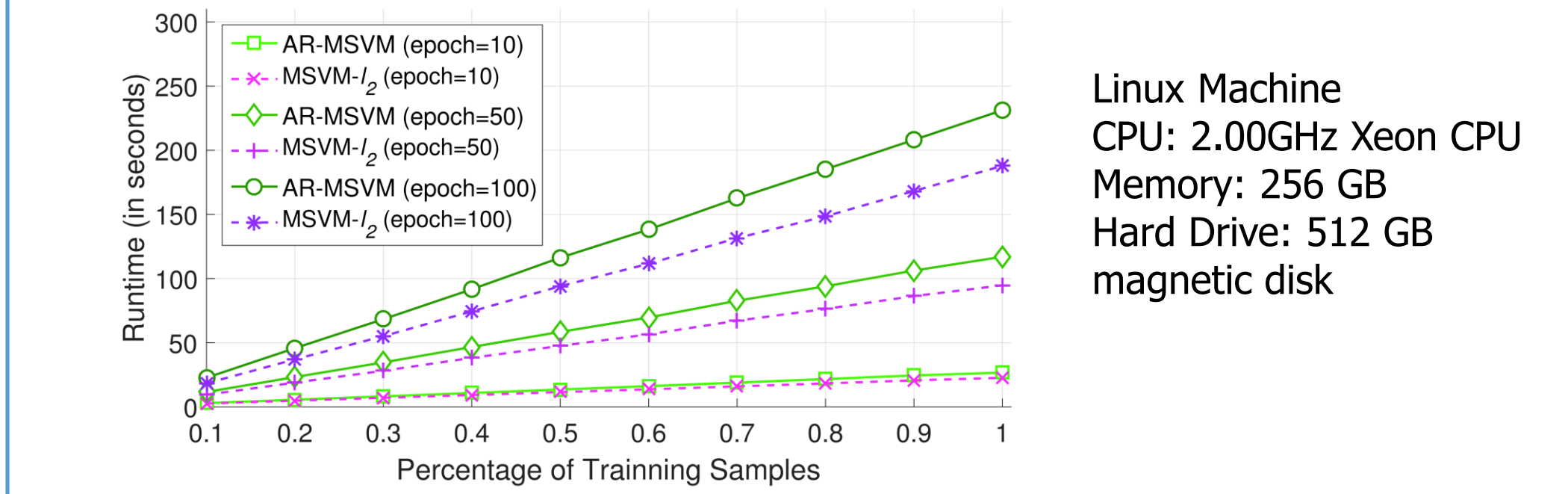
Angular regularizer can improve the performance of MSVM.

- Ensemble methods: efficiently improve the generalization performance.
- L1-regularized methods: sparse constraints are not effective in reduce overfitting of multi-class classification.
- Diversity-promoting methods: CS-MSVM is smooth, while the others are non-smooth. AR-MSVM is not only smooth but also ensures all angles are increased directly.
- L2-regularized methods: AR-MSVM is directly based on MSVM-L2, but achieves better results in all cases.

Parameter Sensitivity:



Runtime:



Imbalanced Classification

Table 4: Classification accuracy (%) on four imbalance datasets

Dataset	Metric	Global-CS	CS	Static-SMT	SL-SMT	AdaB-NC	MSVM- $\ell_2$	AR-MSVM	$\Delta$
Dermatology	Acc	95.78	95.44	95.60	94.30	97.08	84.29 ± 8.99	94.57 ± 1.62	+10.28
Page-Blocks	Acc	91.67	89.32	69.04	89.34	88.29	88.93 ± 2.27	92.27 ± 1.09	+3.34
Wine-Quality-Red	Acc	39.33	37.82	30.74	37.93	37.22	49.94 ± 5.06	52.94 ± 1.79	+3.00
Zoo	Acc	95.02	93.02	95.35	93.02	95.02	93.28 ± 2.00	95.38 ± 1.50	+2.10

A benchmark established by FernNdez et al, 2013

- On small classes (2-5), AR-MSVM performs much better than MSVM.
- On large class (1), they are comparable.

## VI. THEORETICAL RESULTS

We analyze the generalization error under a supervised learning framework via the PAC learning theory.

**Theorem 1.** Fix  $p > 0$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , the following multi-class classification generalization bound holds for all  $h \in H \rightsquigarrow H^1$ :

$$L(h) - L(h^*) \leq \frac{8K^2 C_1 C_2}{p\sqrt{m}} + \frac{2}{\sqrt{m}} + \left(\frac{1}{p}\sqrt{1 - \frac{1}{K}}\right) (\cos(-\mathcal{R}(\mathbf{W})) + \frac{K+1}{K-1} C_1 C_2 + 1) \sqrt{\frac{2 \log(2/\delta)}{m}}.$$

Main conclusions:

- The MSVM's objective function is seeking the maximum margin principle which encourages  $H$  approaching  $H^1$ .
- The generalization error bound  $L(h) - L(h^*)$  is an increasing function of the angular regularizer  $\mathcal{R}(\mathbf{W})$ . It is also why  $\beta$  needs to be chosen to achieve the best accuracy.
- Thus, decreasing angular regularizer  $\mathcal{R}(\mathbf{W})$  can reduce the error bound and improve the generalization performance.

## ACKNOWLEDGEMENTS

This work is supported by China 973 Fundamental R&D Program (No.2014CB340300), NSFC program (No.61472022, 61421003), SKLSD-2016ZX-11, and partly by the Beijing Advanced Innovation Center for Big Data and Brain Computing.

Appendix Files

Wechat Contact

Homepage

