# Multi-Class Learning: From Theory to Algorithm

Jian Li[1,2], **Yong Liu**[1,*], Rong Yin[1,2], Hua Zhang[1], Lizhong Ding[5], Weiping Wang[1,3,4]

[1]Institute of Information Engineering, Chinese Academy of Sciences, China
[2]School of Cyber Security, University of Chinese Academy of Sciences, China
[3]National Engineering Research Center for Information Security
[4]National Engineering Laboratory for Information Security Technology
[5]Inception Institute of Artificial Intelligence (IIAI), Abu Dhabi, UAE

中国科学院 信息工程研究所
INSTITUTE OF INFORMATION ENGINEERING,CAS

## Introduction

▶ Statistical learning of multi-class classification is a crucial problem in machine learning.

▶ Existing generalization bounds for multi-class classification:

| Methods | Convergence rate |
|---|---|
| VC-dimension | $\mathcal{O}(\sqrt{V}\log K/\sqrt{n})$ |
| Natarajan dimension | $\mathcal{O}(d_{Nat}/n)$ |
| Covering Number | $\mathcal{O}(1/\sqrt{n})$ |
| Rademacher Complexity | $\mathcal{O}(\log^2 K/\sqrt{n})$ |
| Stability | $\mathcal{O}(1/\sqrt{n})$ |
| PAC-Bayesian | $\mathcal{O}(\sqrt{\hat{L}(h_\gamma)/n})$ |

▶ Contributions:
  ▷ A new local Rademacher complexity based bound with fast convergence rate $\mathcal{O}\big((\log K)^{2+1/\log K}/n\big)$ for multi-class classification is establish.
  ▷ Two novel multi-class multiple kernel learning algorithms are proposed with statistical guarantees: a) Conv-MKL. b) SMSD-MKL.

## Notations and Preliminaries

▶ Multi-class classification setting
Let $\mathcal{X}$ be the input space and $\mathcal{Y} = \{1, 2, \ldots, K\}$ the output space. Based on training examples $\mathcal{S}$ drawn i.i.d. from a fixed, but unknown probability distribution on $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, we wish to learn a scoring rule $h$ mapping from $\mathcal{Z}$ to $\mathbb{R}$ to predict $x \to \arg\max_{y \in \mathcal{Y}} h(x, y)$. For any $h \in \mathcal{H}$, the margin of a labeled example $z = (x, y)$ is defined as

$$\rho_h(z) := h(x, y) - \max_{y' \neq y} h(x, y').$$

The $h$ misclassifies the labeled example $z = (x, y)$ if $\rho_h(z) \leq 0$. Let $\ell(\rho_h(z))$ be loss function, $L(\ell_h)$ and $\hat{L}(\ell_h)$ be expected generalization error and empirical error with respect to $\ell_h$

$$L(\ell_h) := \mathbb{E}_\mu[\ell(\rho_h(z))] \text{ and } \hat{L}(\ell_h) = \frac{1}{n}\sum_{i=1}^n \ell(\rho_h(z_i)).$$

▶ Hypothesis Space
Let $\kappa : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ be a Mercer kernel with $\phi$ being the associated feature map. The $\ell_p$-norm hypothesis space is denoted by:

$$\mathcal{H}_{p,\kappa} = \Big\{h_w = (\langle w_1, \phi(x)\rangle, \ldots, \langle w_K, \phi(x)\rangle) :$$
$$\|w\|_{2,p} \leq 1, 1 \leq p \leq 2\Big\},$$

where $w = (w_1, \ldots, w_K)$ and $\|w\|_{2,p} = \left[\sum_{i=1}^K \|w_i\|_2^p\right]^{1/p}$ is the $\ell_{2,p}$-norm. For $p \geq 1$, the dual exponent $q$ satisfies $1/p + 1/q = 1$. The space of loss function associated with $\mathcal{H}_{p,\kappa}$ is denoted by

$$\mathcal{L} = \{\ell_h := \ell(\rho_h(z)) : h \in \mathcal{H}_{p,\kappa}\}.$$

▶ The local Rademacher complexity of $\mathcal{L}$

$$\mathcal{R}(\mathcal{L}^r) := \mathcal{R}\left\{a\ell_h \Big| a \in [0, 1], \ell_h \in \mathcal{L}, L[(a\ell_h)^2] \leq r\right\}.$$

## Sharper Generalization Bounds

▶ **Local Rademacher complexity of multi-class classification**
With probability at least $1 - \delta$,

$$\mathcal{R}(\mathcal{L}^r) \leq \frac{c_{d,\vartheta}\xi(K)\sqrt{\zeta r}\log^{\frac{3}{2}}(n)}{\sqrt{n}} + \frac{4\log(1/\delta)}{n},$$

where

$$\xi(K) = \begin{cases} \sqrt{e}(4\log K)^{1+\frac{1}{2\log K}}, & \text{if } q \geq 2\log K, \\ (2q)^{1+\frac{1}{q}}K^{\frac{1}{q}}, & \text{otherwise,} \end{cases}$$

$c_{d,\vartheta}$ is a constant depending on $d$ and $\vartheta$.

▶ **A Sharper Generalization Bound**
$\forall h \in \mathcal{H}_{p,\kappa}$ and $\forall k > \max(1, \frac{\sqrt{2}}{2d})$, with probability at least $1 - \delta$,

$$L(h) \leq \max\left\{\frac{k}{k-1}\hat{L}(\ell_h), \hat{L}(\ell_h) + \frac{c_{d,\vartheta,\zeta,k}\xi^2(K)\log^3 n}{n} + \frac{c_\delta}{n}\right\},$$

where

$$\xi(K) = \begin{cases} \sqrt{e}(4\log K)^{1+\frac{1}{2\log K}}, & \text{if } q \geq 2\log K, \\ (2q)^{1+\frac{1}{q}}K^{\frac{1}{q}}, & \text{otherwise,} \end{cases}$$

const $c_{d,\vartheta}$ depends on $d, \vartheta, \zeta, k$, and const $c_\delta$ depends $\delta$.

## Multi-Class Multiple Kernel Learning

▶ **Conv-MKL**
Consider use multiple kernels $\kappa_\mu = \sum_{m=1}^M \mu_m \kappa_m$, the $\ell_p$ hypothesis space of multiple kernels can be written as:

$$\mathcal{H}_{mkl} = \Big\{h_{w,\kappa_\mu} = (\langle w_1, \phi_\mu(x)\rangle, \ldots, \langle w_K, \phi_\mu(x)\rangle),$$
$$\|w\|_{2,p} \leq 1, 1 < p \leq \frac{2\log K}{2\log K - 1}\Big\}.$$

According to theoretical analysis, we add local Rademacher complexity (the tail sum of the eigenvalues of the kernel) to restrict $\mathcal{H}_{mkl}$:

$$\mathcal{H}_2 = \Big\{h_{w,\kappa_\mu} \in \mathcal{H}_{mkl} : \sum_{m=1}^M \mu_m \sum_{j>\zeta} \lambda_j(K_m) \leq 1\Big\}.$$

Using normalized kernels $\tilde{\kappa}_m = \left(\sum_{j>\zeta} \lambda_j(K_m)\right)^{-1} \kappa_m$ and $\tilde{\kappa}_\mu = \sum_{m=1}^M \mu_m \tilde{\kappa}_m$, we can simply rewrite $\mathcal{H}_2$ as

$$\Big\{h_{w,\tilde{\kappa}_\mu} = (\langle w_1, \tilde{\phi}_\mu(x)\rangle, \ldots, \langle w_K, \tilde{\phi}_\mu(x)\rangle),$$
$$\|w\|_{2,p} \leq 1, 1 < p \leq \frac{2\log K}{2\log K - 1}, \mu \succeq 0, \|\mu\|_1 \leq 1\Big\},$$

With precomputed kernel matrices regularized by local Rademacher complexity, the method gets solution by any $\ell_p$-norm MC-MKL solvers.

▶ **SMSD-MKL**
Considering a more challenging case, we perform penalized ERM over the class $\mathcal{H}_1$, aiming to solve a convex optimization problem with an additional term representing local Rademacher complexity :

$$\min_{w,\mu} \underbrace{\frac{1}{n}\sum_{i=1}^n \ell(w, \phi_\mu(x_i), y_i)}_{C(w)} + \underbrace{\frac{\alpha}{2}\|w\|_{2,p}^2 + \beta\sum_{m=1}^M \mu_m r_m}_{\Omega(w)},$$

where
$$\ell(w, \phi_\mu(x_i), y_i) = \left|1 - \left(\langle w_{y_i}, \phi_\mu(x_i)\rangle - \max_{y \neq y_i}\langle w_y, \phi_\mu(x_i)\rangle\right)\right|_+ \text{ and }$$
$r_m = \sum_{j>\zeta} \lambda_j(K_m)$ is the tail sum of the eigenvalues of the $m$-th kernel matrix, $m = 1, \ldots, M$.

Based on widely used stochastic mirror descent framework, we design SMSD-MKL algorithm, implemented by stochastic sub-gradient descent with updating dual weights, to solve above optimization objective.

## Experiments

Table: Comparison of average test accuracies of our Conv-MKL and SMSD-MKL with the others. We bold the numbers of the best method and underline the numbers of the other methods which are not significantly worse than the best one.

| | Conv-MKL | SMSD-MKL | LMC | One vs. One | One vs. Rest | GMNP | $\ell_1$ MC-MKL | $\ell_2$ MC-MKL | UFO-MKL |
|---|---|---|---|---|---|---|---|---|---|
| plant | 77.14±2.25 | **78.01±2.17** | 70.12±2.96 | 75.83±2.69 | 75.17±2.68 | 75.42±3.64 | 77.60±2.63 | 75.49±2.48 | 76.77±2.42 |
| psortPos | 74.41±3.35 | **76.23±3.39** | 63.85±3.94 | 73.33±4.21 | 71.70±4.89 | 73.55±4.22 | 71.87±4.87 | 70.70±4.89 | 74.56±4.04 |
| psortNeg | 74.07±2.16 | **74.66±1.90** | 57.85±2.49 | 73.74±2.87 | 71.94±2.50 | 74.27±2.51 | 72.83±2.20 | 72.42±2.65 | 73.80±2.26 |
| nonpl | **79.15±1.51** | 78.69±1.58 | 75.16±1.48 | 77.78±1.52 | 77.49±1.53 | 78.35±1.46 | 77.89±1.79 | 77.95±1.64 | 78.07±1.56 |
| sector | 92.83±2.62 | **93.39±0.70** | 93.16±0.66 | 90.61±0.69 | 91.34±0.61 | \ | \ | 92.15±2.57 | 92.60±0.47 |
| segment | 96.79±0.91 | **97.62±0.83** | 95.07±1.11 | 97.08±0.61 | 97.02±0.80 | 96.87±0.80 | 96.98±0.64 | 97.58±0.68 | 97.20±0.82 |
| vehicle | **79.35±2.27** | 77.28±2.78 | 75.61±3.56 | 78.72±1.92 | 79.11±1.94 | 81.57±2.24 | 74.96±2.93 | 76.27±3.15 | 76.92±2.83 |
| vowel | 98.82±1.19 | **98.83±5.57** | 62.32±4.97 | 98.12±1.76 | 98.22±1.83 | 97.04±1.85 | 98.27±1.22 | 97.86±1.75 | 98.22±1.62 |
| wine | **99.63±0.96** | **99.63±0.96** | 97.87±2.80 | 97.24±3.05 | 98.14±3.04 | 97.69±2.43 | 98.61±1.75 | 98.52±1.89 | 99.44±1.13 |
| dna | 96.08±0.83 | **96.30±0.79** | 92.02±1.50 | 95.89±0.56 | 95.61±0.73 | 94.60±0.94 | 96.27±0.68 | 95.06±0.92 | 95.84±0.61 |
| glass | **75.19±5.05** | 73.72±5.80 | 63.95±6.04 | 71.98±5.75 | 70.00±5.75 | 71.24±8.14 | 69.07±8.08 | 74.03±6.41 | 72.46±6.12 |
| iris | 96.67±2.94 | **97.00±2.63** | 88.00±7.82 | 95.93±3.25 | 95.87±3.20 | 95.40±7.34 | 95.40±6.46 | 94.00±7.82 | 95.93±2.88 |
| svmguide2 | 82.69±5.65 | **85.17±3.83** | 81.10±4.15 | 84.79±3.45 | 84.27±3.03 | 81.77±3.45 | 83.16±3.63 | 83.84±4.21 | 82.91±3.09 |
| satimage | 91.64±0.88 | 91.78±0.82 | 84.95±1.15 | 90.67±0.91 | 89.29±0.96 | 89.97±0.81 | 91.86±0.62 | 90.43±1.27 | **91.92±0.83** |

## Conclusions

▶ A new local Rademacher complexity based bound with fast convergence rate for multi-class classification is establish. Convergence rate is improved from sub-linear to linear $\mathcal{O}(\sqrt{\frac{K}{n}}) \Rightarrow \mathcal{O}(\frac{(\log K)^{2+1/\log K}}{n})$.

▶ Two novel multi-class classification algorithms are proposed with statistical guarantees: a) Conv-MKL. b) SMSD-MKL.

## Main References

[1] P. L. Bartlett, O. Bousquet, and S. Mendelson.
Local Rademacher complexities.
The Annals of Statistics, 33(4):1497–1537, 2005.

[2] F. Orabona and J. Luo.
Ultra-fast optimization algorithm for sparse multi kernel learning.
Proceedings of the 28th International Conference on Machine Learning (ICML 2011), pages 249–256, 2011.