

Appendix for “Can Infinitely Wide Deep Nets Help Small-data Multi-label Learning?”

A Proof of Proposition 1

Proposition 1. Consider minimizing $C(\boldsymbol{\theta})$ in Eq.(3) by gradient descent with infinitesimally small learning rate (i.e., gradient flow): $\frac{d\boldsymbol{\theta}(t)}{dt} = -\nabla C(\boldsymbol{\theta}(t))$. Let $\mathbf{F}(t) = [f(\mathbf{x}_1, \boldsymbol{\theta}(t)); \dots; f(\mathbf{x}_n, \boldsymbol{\theta}(t))] \in \mathbb{R}^{n \times c}$ be network outputs on all \mathbf{x}_i 's at time t . When the hidden widths $d_1, \dots, d_L \rightarrow \infty$, then $\mathbf{F}(t)$ follows the following evolution:

$$\frac{d\mathbf{F}(t)}{dt} = -\frac{1}{n} \mathbf{K} \cdot \frac{\partial \ell(t)}{\partial \mathbf{F}(t)}, \quad (1)$$

where $\mathbf{K} \in \mathbb{R}^{n \times n}$ is the NTK kernel matrix and $\frac{\partial \ell(t)}{\partial \mathbf{F}(t)} \triangleq \left[\frac{\partial \ell(f(\mathbf{x}_1, \boldsymbol{\theta}(t)), \mathbf{y}_1)}{\partial f(\mathbf{x}_1, \boldsymbol{\theta}(t))}; \dots; \frac{\partial \ell(f(\mathbf{x}_n, \boldsymbol{\theta}(t)), \mathbf{y}_n)}{\partial f(\mathbf{x}_n, \boldsymbol{\theta}(t))} \right] \in \mathbb{R}^{n \times c}$.

Proof. First, when discarding the infinity width assumption, we can get the dynamics of $\mathbf{F}(t)$ w.r.t. a matrix-based kernel matrix $\Theta^{(L)}(\boldsymbol{\theta}) \in \mathbb{R}^{nc \times nc}$, where its corresponding matrix-based kernel function is $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^{c \times c}$. Besides, the kernel matrix $\Theta^{(L)}(\boldsymbol{\theta})$ depends on the parameters $\boldsymbol{\theta}$ and varies during training.

Second, when adding the infinity width assumption, based on the Theorem 1 and 2 in [4], we can get the deterministic kernel matrix $\Theta^{(L)} = \mathbf{K} \otimes \mathbf{I}_c$, where $\mathbf{I}_c \in \mathbb{R}^{c \times c}$ is the identity matrix. The kernel matrix $\Theta^{(L)}$ doesn't vary during training and always equals the state in the random initialization. Besides, the matrix-based kernel is separable and can be equivalently transformed into the corresponding scalar-based kernel [1], which is used here. \square

B Derivation of the gradient $\nabla g_i(\mathbf{W})$

Here, we take the logistic base loss function for example, which is used in our experiments.

Define

$$g_i(\mathbf{W}) = \frac{1}{c} \sum_{j=1}^c \log(1 + \exp(-y_{ij} \langle \mathbf{w}^j, z(\mathbf{x}_i) \rangle)) + \frac{\tau}{2} \|\mathbf{W}\|_F^2 + \bar{g}_i(\mathbf{W}), \quad (2)$$

where

$$\bar{g}_i(\mathbf{W}) = \frac{\lambda}{|Y_i^+| |Y_i^-|} \sum_{p \in Y_i^+} \sum_{q \in Y_i^-} \log(1 + \exp(\langle \mathbf{w}^q - \mathbf{w}^p, z(\mathbf{x}_i) \rangle)). \quad (3)$$

Then, the gradient of $g_i(\mathbf{W})$ w.r.t. \mathbf{W} is as follows:

$$\nabla g_i(\mathbf{W}) = \frac{1}{c} \mathbf{x}_i^\top \left(\frac{1}{\mathbf{1} + \exp(-\mathbf{y}_i \circ (z(\mathbf{x}_i) \mathbf{W}))} \circ \exp(-\mathbf{y}_i \circ (z(\mathbf{x}_i) \mathbf{W})) \right) + \tau \mathbf{W} + \nabla \bar{g}_i(\mathbf{W}), \quad (4)$$

where $\nabla \bar{g}_i(\mathbf{W}) = \left[\frac{\partial \bar{g}_i}{\partial \mathbf{w}^1}, \dots, \frac{\partial \bar{g}_i}{\partial \mathbf{w}^c} \right]$ and for $j \in [c]$,

$$\frac{\partial \bar{g}_i}{\partial \mathbf{w}^j} = \frac{\lambda}{|Y_i^+| |Y_i^-|} \left\{ \mathbb{1}[j \in Y_i^+] \sum_{q \in Y_i^-} \frac{-z(\mathbf{x}_i)}{1 + \exp(\langle \mathbf{w}^j - \mathbf{w}^q, z(\mathbf{x}_i) \rangle)} + \mathbb{1}[j \in Y_i^-] \sum_{p \in Y_i^+} \frac{z(\mathbf{x}_i)}{1 + \exp(\langle \mathbf{w}^p - \mathbf{w}^j, z(\mathbf{x}_i) \rangle)} \right\}. \quad (5)$$

C Proofs for the generalization analyses

Technically, we provide the generalization analyses based on Rademacher complexity [2] and the recent vector-contraction inequality [5], following recent work [6].

First, we define the surrogate expected loss w.r.t. Hamming Loss (HL) and Ranking Loss (RL) as follows.

$$R^h(f) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P} \left[L_h(f(\mathbf{x}), \mathbf{y}_i) \right], \quad R^r(f) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P} \left[L_r(f(\mathbf{x}), \mathbf{y}_i) \right]. \quad (6)$$

Then, we give the following for the subsequent analyses.

Theorem C.1 (The base theorem [6]). *Assume the loss function $L : \mathbb{R}^c \times \{-1, +1\}^c \rightarrow \mathbb{R}_+$ is μ -Lipschitz continuous w.r.t. the first argument and bounded by M . Then, for any $\delta > 0$, with probability at least $1 - \delta$ over the draw of an i.i.d. sample S of size n , the following generalization bound holds for all $f \in \mathcal{F}$:*

$$R(f) \leq \hat{R}_S(f) + 2\sqrt{2}\mu\hat{\mathfrak{R}}_S(\mathcal{F}) + 3M\sqrt{\frac{\log \frac{2}{\delta}}{2n}}. \quad (7)$$

Lemma C.1 (The Rademacher complexity of the kernel-based hypothesis set). *Consider the hypothesis set $\mathcal{F} = \{ \langle f_1^\kappa, \Phi(\mathbf{x}) \rangle, \dots, \langle f_c^\kappa, \Phi(\mathbf{x}) \rangle : \sum_{j=1}^c \|f_j^\kappa\|_{\mathbb{H}}^2 \leq \Lambda \}$, which can be equivalently transformed into the following form $\mathcal{F} = \{ \mathbf{x} \mapsto \mathbf{W}^\top \phi(\mathbf{x}) : \mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_c)^\top, \|\mathbf{W}\|_{\mathbb{H}, 2} \leq \Lambda \}$, $\hat{\mathfrak{R}}_S(\mathcal{F})$, where $\|\mathbf{W}\|$ denotes $\|\mathbf{W}\|_{\mathbb{H}, 2}$ for convenience. Then, $\hat{\mathfrak{R}}_S(\mathcal{F})$ can be bounded below:*

$$\hat{\mathfrak{R}}_S(\mathcal{F}) \leq \sqrt{\frac{c \max\{\mathbf{K}_{ii}\} \Lambda^2}{n}}. \quad (8)$$

Proof. For the kernel-based hypothesis set $\mathcal{F} = \{ \mathbf{x} \mapsto \mathbf{W}^\top \phi(\mathbf{x}) : \mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_c)^\top, \|\mathbf{W}\| \leq \Lambda \}$, the following inequalities about $\hat{\mathfrak{R}}_S(\mathcal{F})$ hold:

$$\begin{aligned} \hat{\mathfrak{R}}_S(\mathcal{F}) &= \frac{1}{n} \mathbb{E}_\epsilon \left[\sup_{\|\mathbf{W}\| \leq \Lambda} \sum_{i=1}^n \sum_{j=1}^c \epsilon_{ij} \langle \mathbf{w}_j, \phi(\mathbf{x}_i) \rangle \right] \\ &= \frac{1}{n} \mathbb{E}_\epsilon \left[\sup_{\|\mathbf{W}\| \leq \Lambda} \sum_{j=1}^c \langle \mathbf{w}_j, \sum_{i=1}^n \epsilon_{ij} \phi(\mathbf{x}_i) \rangle \right] \\ &= \frac{1}{n} \mathbb{E}_\epsilon \left[\sup_{\|\mathbf{W}\| \leq \Lambda} \langle \mathbf{W}, \mathbf{X}_\epsilon \rangle \right] \quad (\mathbf{X}_\epsilon = [\sum_{i=1}^n \epsilon_{i1} \phi(\mathbf{x}_i), \dots, \sum_{i=1}^n \epsilon_{ic} \phi(\mathbf{x}_i)]) \\ &\leq \frac{1}{n} \mathbb{E}_\epsilon \left[\sup_{\|\mathbf{W}\| \leq \Lambda} \|\mathbf{W}\| \|\mathbf{X}_\epsilon\| \right] \quad (\text{Cauchy-Schwarz Inequality}) \\ &= \frac{\Lambda}{n} \mathbb{E}_\epsilon \left[\sum_{j=1}^c \left\| \sum_{i=1}^n \epsilon_{ij} \phi(\mathbf{x}_i) \right\|^2 \right]^{1/2} \\ &= \frac{\Lambda}{n} \mathbb{E}_\epsilon \left[\sum_{j=1}^c \sum_{p=1}^n \sum_{q=1}^n \epsilon_{pj} \epsilon_{qj} \langle \phi(\mathbf{x}_p), \phi(\mathbf{x}_q) \rangle \right]^{1/2} \\ &= \frac{\Lambda}{n} \mathbb{E}_\epsilon \left[\sum_{j=1}^c \sum_{i=1}^n \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_i) \rangle \right]^{1/2} \quad (\forall p \neq q, \mathbb{E}[\epsilon_{pj} \epsilon_{qj}] = \mathbb{E}[\epsilon_{pj}] \mathbb{E}[\epsilon_{qj}] = 0 \text{ and } \mathbb{E}[\epsilon_{ij} \epsilon_{ij}] = 1) \\ &= \frac{\Lambda \sqrt{c \operatorname{Tr}(\mathbf{K})}}{n} \quad (\kappa(\mathbf{x}_i, \mathbf{x}_i) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_i) \rangle, \mathbf{K} = [\kappa(\mathbf{x}_i, \mathbf{x}_j)] \text{ is the kernel matrix}) \\ &\leq \sqrt{\frac{c \max\{\mathbf{K}_{ii}\} \Lambda^2}{n}}. \end{aligned} \quad (9)$$

□

Lemma C.2 (The property of the surrogate loss function [6]). *Assume that the base loss function $\ell(u)$ is ρ -Lipschitz continuous and bounded by B . Then, the surrogate Hamming Loss L_h is $\frac{\rho}{\sqrt{c}}$ -Lipschitz, and the surrogate Ranking Loss L_r is ρ -Lipschitz w.r.t. the first argument. Besides, they are all bounded by B .*

C.1 Proof of Theorem 5

Theorem 5 (Learning guarantee w.r.t. the Hamming Loss). *Consider the hypothesis space $\mathcal{F}_h = \{ f \in \mathcal{F} : \hat{R}_S^r(f) \leq \Lambda_1 \}$ and the loss function L_h . Besides, Assumption 1 holds. Then, for any $\delta > 0$, the following holds for any $f \in \mathcal{F}_h$ with*

probability at least $1 - \delta$

$$R_{0/1}^h(f) \leq \hat{R}_S^h(f) + 2\rho\sqrt{\frac{2}{c}}\hat{\mathfrak{R}}_S(\mathcal{F}_h) + 3B\sqrt{\frac{\log \frac{2}{\delta}}{2n}}, \quad (10)$$

where the empirical Rademacher complexity satisfies

$$\hat{\mathfrak{R}}_S(\mathcal{F}_h) \leq \hat{\mathfrak{R}}_S(\mathcal{F}) \leq \sqrt{\frac{c \max\{\mathbf{K}_{ii}\}\Lambda^2}{n}}. \quad (11)$$

Proof. Since $L = L_h$, we can get its Lipschitz constant (i.e. $\frac{\rho}{\sqrt{c}}$) and bounded value (i.e. B) from Lemma C.2. Then, applying Theorem C.1 and the inequality $R_{0/1}^h(\text{sgn} \circ f) \leq R^h(f)$, we can Eq.(10). Besides, it can be easily verified that $\hat{\mathfrak{R}}_S(\mathcal{F}_h) \leq \hat{\mathfrak{R}}_S(\mathcal{F})$. Then, based on Lemma C.1, we can get Eq.(11). \square

C.2 Proof of Theorem 6

Theorem 6 (Learning guarantee w.r.t. the Ranking Loss). *Consider the hypothesis space $\mathcal{F}_r = \{f \in \mathcal{F} : \hat{R}_S^h(f) \leq \Lambda_2\}$ and the loss function L_r . Besides, Assumption 1 holds. Then, for any $\delta > 0$, the following holds for any $f \in \mathcal{F}_r$ with probability at least $1 - \delta$*

$$R_{0/1}^r(f) \leq \hat{R}_S^r(f) + 2\sqrt{2}\rho\hat{\mathfrak{R}}_S(\mathcal{F}_r) + 3B\sqrt{\frac{\log \frac{2}{\delta}}{2n}}, \quad (12)$$

where the empirical Rademacher complexity satisfies

$$\hat{\mathfrak{R}}_S(\mathcal{F}_r) \leq \hat{\mathfrak{R}}_S(\mathcal{F}) \leq \sqrt{\frac{c \max\{\mathbf{K}_{ii}\}\Lambda^2}{n}}. \quad (13)$$

Proof. Since $L = L_r$, we can get its Lipschitz constant (i.e. ρ) and bounded value (i.e. B) from Lemma C.2. Then, applying Theorem C.1 and the inequality $R_{0/1}^r(\text{sgn} \circ f) \leq R^h(f)$, we can Eq.(12). Besides, it can be easily verified that $\hat{\mathfrak{R}}_S(\mathcal{F}_h) \leq \hat{\mathfrak{R}}_S(\mathcal{F})$. Then, based on Lemma C.1, we can get Eq.(13). \square

D Additional experimental results

D.1 Results w.r.t. other measures

The experimental results w.r.t. the measures of Subset Accuracy, instance-F1, Coverage, and Average Precision are summarized in Table 1, 2, 3 and 4, respectively.

Table 1: Experimental results of benchmark approaches (mean_{std}) w.r.t. Subset Accuracy (\uparrow) on all datasets. Best results are highlighted in bold.

Dataset	Rank-SVM	BP-MLL	BR-SVM	CPNL	MLFE	MLC-RBF	MLC-NTK
flags	0.073 _{0.025}	0.107 _{0.015}	0.132 _{0.041}	0.156_{0.053}	0.151 _{0.041}	0.144 _{0.035}	0.153 _{0.031}
birds	0.148 _{0.018}	0.121 _{0.019}	0.214 _{0.028}	0.214 _{0.026}	0.170 _{0.018}	0.206 _{0.024}	0.234_{0.025}
emotions	0.291 _{0.025}	0.268 _{0.039}	0.313 _{0.015}	0.324_{0.035}	0.291 _{0.035}	0.324_{0.025}	0.310 _{0.027}
image	0.451 _{0.018}	0.217 _{0.024}	0.482 _{0.018}	0.533 _{0.017}	0.463 _{0.015}	0.539 _{0.008}	0.563_{0.021}
scene	0.563 _{0.018}	0.587 _{0.019}	0.655 _{0.009}	0.699 _{0.014}	0.617 _{0.009}	0.713 _{0.010}	0.735_{0.009}
yeast	0.156 _{0.012}	0.111 _{0.007}	0.190 _{0.009}	0.179 _{0.006}	0.172 _{0.014}	0.176 _{0.007}	0.202_{0.009}
enron	0.119 _{0.033}	0.099 _{0.008}	0.128 _{0.013}	0.128 _{0.006}	0.124 _{0.014}	0.100 _{0.009}	0.142_{0.010}
business	0.454 _{0.085}	0.359 _{0.003}	0.565 _{0.010}	0.557 _{0.008}	0.538 _{0.006}	0.542 _{0.010}	0.569_{0.008}

D.2 Computational costs

The computational costs of the six algorithms on benchmark datasets are illustrated in Fig. 1. Note that the CPU time is plotted in the log scale in Figure 1. Besides, BR-SVM is implemented that based on the efficient LibSVM library [3] and thus is more efficient than the others implemented by the original Matlab code. For the last datasets, we use the linear kernel for CPNL, where it has more efficient implementation that is tailored for the linear model. In comparison, for the first six datasets, Rank-SVM and CPNL involves the kernel matrix in training. In this case, we can observe that when the number of the dataset becomes relatively larger (i.e., Image, scene, and yeast), our method MLC-NTK (with Nyström method) is more efficient than these two. Note that in experiments we set $m = r = n$ in the Nyström method for high performance although it can reduce r (and m) to accelerate the efficiency.

Table 2: Experimental results of benchmark approaches (mean_{std}) w.r.t. instance-F1 (\uparrow) on all datasets. Best results are highlighted in bold.

Dataset	Rank-SVM	BP-MLL	BR-SVM	CPNL	MLFE	MLC-RBF	MLC-NTK
flags	0.703 _{0.017}	0.717_{0.004}	0.700 _{0.017}	0.697 _{0.031}	0.688 _{0.011}	0.703 _{0.017}	0.711 _{0.018}
birds	0.338 _{0.023}	0.425 _{0.014}	0.469 _{0.017}	0.476_{0.031}	0.400 _{0.029}	0.417 _{0.030}	0.468 _{0.024}
emotions	0.645 _{0.020}	0.655 _{0.034}	0.620 _{0.020}	0.684_{0.019}	0.621 _{0.021}	0.656 _{0.022}	0.656 _{0.026}
image	0.631 _{0.016}	0.523 _{0.026}	0.623 _{0.014}	0.698 _{0.012}	0.593 _{0.014}	0.682 _{0.010}	0.700_{0.013}
scene	0.664 _{0.015}	0.730 _{0.011}	0.717 _{0.010}	0.802_{0.009}	0.685 _{0.010}	0.783 _{0.006}	0.795 _{0.009}
yeast	0.632_{0.007}	0.614 _{0.015}	0.623 _{0.006}	0.630 _{0.007}	0.607 _{0.011}	0.614 _{0.007}	0.629 _{0.005}
enron	0.563 _{0.026}	0.566 _{0.006}	0.529 _{0.010}	0.585_{0.007}	0.538 _{0.012}	0.491 _{0.007}	0.566 _{0.008}
business	0.734 _{0.036}	0.722 _{0.001}	0.763 _{0.005}	0.770_{0.006}	0.765 _{0.005}	0.739 _{0.006}	0.770_{0.005}

Table 3: Experimental results of benchmark approaches (mean_{std}) w.r.t. Coverage (\downarrow) on all datasets. Best results are highlighted in bold.

Dataset	Rank-SVM	BP-MLL	BR-SVM	CPNL	MLFE	MLC-RBF	MLC-NTK
flags	0.525 _{0.035}	0.551 _{0.010}	0.548 _{0.019}	0.558 _{0.021}	0.561 _{0.019}	0.531 _{0.021}	0.525_{0.022}
birds	0.256 _{0.014}	0.285 _{0.008}	0.262 _{0.011}	0.254 _{0.019}	0.276 _{0.024}	0.258 _{0.011}	0.235_{0.021}
emotions	0.294 _{0.011}	0.310 _{0.014}	0.386 _{0.017}	0.277 _{0.010}	0.282 _{0.013}	0.279 _{0.008}	0.275_{0.010}
image	0.171 _{0.008}	0.210 _{0.005}	0.227 _{0.008}	0.157_{0.006}	0.165 _{0.006}	0.161 _{0.004}	0.157_{0.007}
scene	0.068 _{0.004}	0.085 _{0.001}	0.119 _{0.004}	0.064 _{0.002}	0.067 _{0.003}	0.065 _{0.002}	0.061_{0.004}
yeast	0.446 _{0.006}	0.480 _{0.010}	0.627 _{0.007}	0.445 _{0.006}	0.452 _{0.006}	0.441 _{0.006}	0.433_{0.005}
enron	0.235 _{0.021}	0.242 _{0.001}	0.580 _{0.012}	0.232 _{0.006}	0.228 _{0.010}	0.245 _{0.005}	0.207_{0.009}
business	0.068 _{0.005}	0.092 _{0.005}	0.338 _{0.009}	0.065 _{0.002}	0.082 _{0.004}	0.069 _{0.003}	0.063_{0.002}

Table 4: Experimental results of benchmark approaches (mean_{std}) w.r.t. Average Precision (\uparrow) on all datasets. Best results are highlighted in bold.

Dataset	Rank-SVM	BP-MLL	BR-SVM	CPNL	MLFE	MLC-RBF	MLC-NTK
flags	0.821 _{0.014}	0.816 _{0.002}	0.808 _{0.016}	0.803 _{0.022}	0.803 _{0.012}	0.824 _{0.019}	0.824_{0.015}
birds	0.635 _{0.015}	0.591 _{0.015}	0.637 _{0.016}	0.650 _{0.020}	0.634 _{0.028}	0.619 _{0.016}	0.656_{0.023}
emotions	0.808 _{0.010}	0.786 _{0.018}	0.760 _{0.015}	0.828 _{0.010}	0.822 _{0.012}	0.827 _{0.010}	0.830_{0.018}
image	0.823 _{0.010}	0.762 _{0.013}	0.772 _{0.011}	0.839 _{0.007}	0.826 _{0.008}	0.832 _{0.007}	0.841_{0.007}
scene	0.882 _{0.008}	0.857 _{0.005}	0.834 _{0.006}	0.893_{0.006}	0.885 _{0.005}	0.885 _{0.004}	0.893_{0.006}
yeast	0.755 _{0.005}	0.732 _{0.022}	0.680 _{0.007}	0.775 _{0.009}	0.769 _{0.008}	0.770 _{0.008}	0.780_{0.005}
enron	0.672 _{0.025}	0.674 _{0.005}	0.482 _{0.010}	0.702 _{0.010}	0.705 _{0.008}	0.663 _{0.005}	0.714_{0.006}
business	0.860 _{0.036}	0.865 _{0.002}	0.742 _{0.005}	0.891 _{0.005}	0.885 _{0.004}	0.878 _{0.004}	0.893_{0.004}

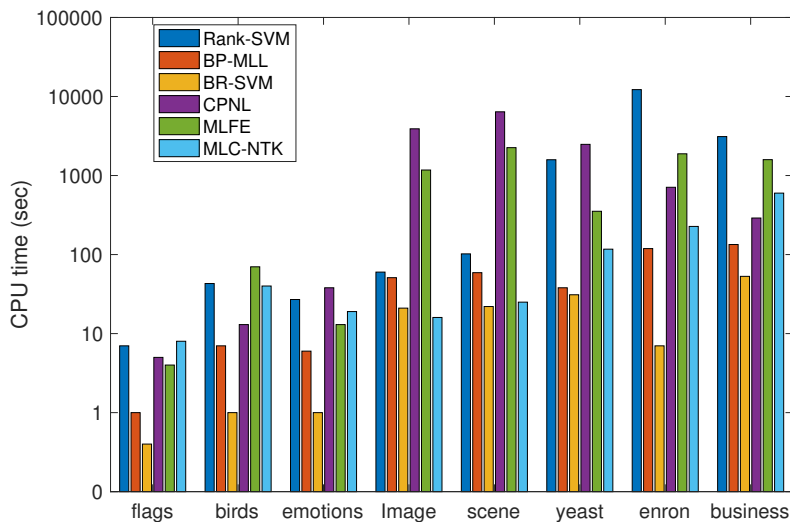


Figure 1: Computational costs of the six algorithms on benchmark datasets.

D.3 Sensitivity analysis w.r.t. the hyper-parameters τ and λ

In this section, we give the sensitivity analysis w.r.t. the hyper-parameters τ and λ , which are illustrated in Fig. 2 and 3 respectively. We can observe that the hyper-parameter τ is more sensitive to λ because we argue that τ has more power to control the model complexity than λ .

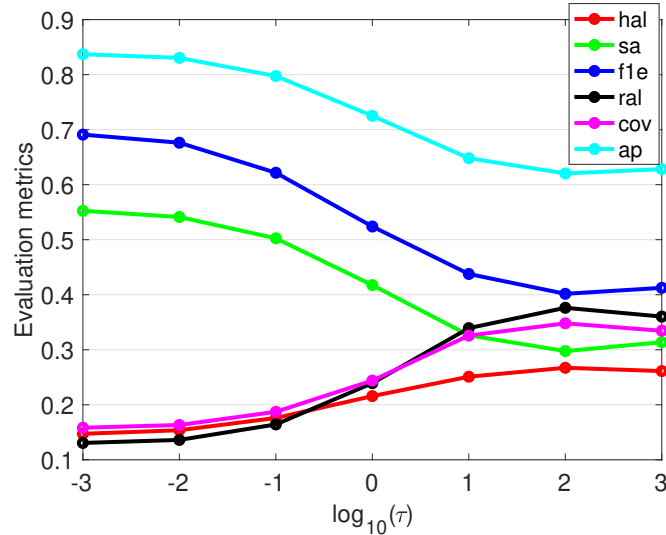


Figure 2: Sensitivity analysis w.r.t. the hyper-parameter τ . Hal is the abbreviation of Hamming Loss, and so on.

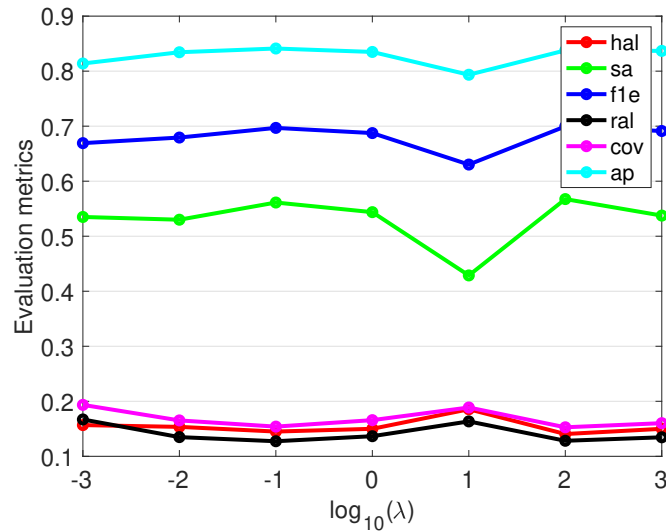


Figure 3: Sensitivity analysis w.r.t. the hyper-parameter λ . Hal is the abbreviation of Hamming Loss, and so on.

References

- [1] Mauricio A Alvarez, Lorenzo Rosasco, Neil D Lawrence, et al. Kernels for vector-valued functions: A review. *Foundations and Trends® in Machine Learning*, 4(3):195–266, 2012.
- [2] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.

- [3] Chih-Chung Chang and Chih-Jen Lin. Libsvm: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):1–27, 2011.
- [4] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018.
- [5] Andreas Maurer. A vector-contraction inequality for rademacher complexities. In *International Conference on Algorithmic Learning Theory*, pages 3–17. Springer, 2016.
- [6] Guoqiang Wu and Jun Zhu. Multi-label classification: do hamming loss and subset accuracy really conflict with each other? *Advances in Neural Information Processing Systems*, 33, 2020.