



ICML

International Conference
On Machine Learning



A General Representation Learning Framework with Generalization Performance Guarantees

Junbiao Cui¹, Jianqing Liang¹, Qin Yue¹, Jiye Liang¹

Speaker: Junbiao Cui

1. Key Laboratory of Computational Intelligence and Chinese Information Processing of Ministry of Education, School of Computer and Information Technology, Shanxi University, Taiyuan 030006, Shanxi, China.

Correspondence to: Jiye Liang <ljy@sxu.edu.cn>



Outline



1. Motivation

2. Proposed Criterion

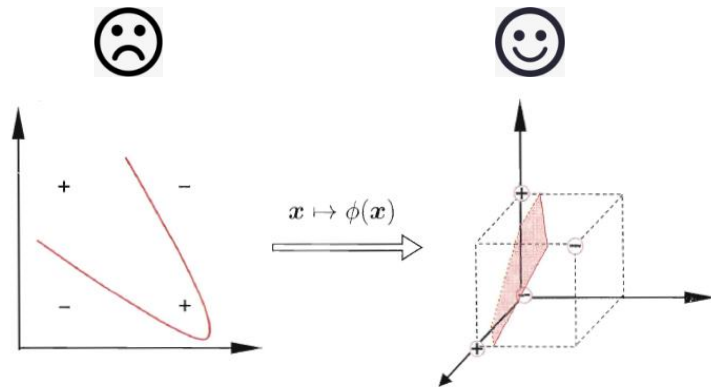
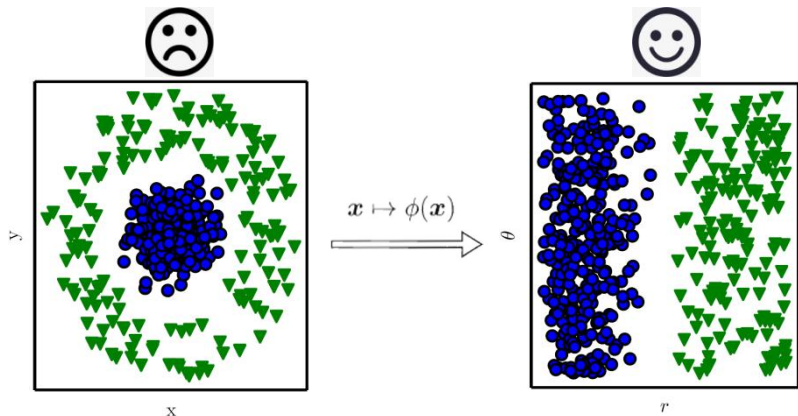
3. Application I: Kernel Selection

4. Application II: DNN Boosting

5. Conclusion and Outlook

1. Motivation

- **Truism** Good data representation leads to good generalization performance



- **However**

General Framework of Machine Learning

$$h^* = \arg \min_{h \in \mathcal{H}} \sum_{i=1}^n \ell(y_i, h(x_i)) + \mathcal{R}(h), \quad h^* : \mathcal{X} \rightarrow \mathcal{Y}$$

The relationship between representation learning and generalization performance is **not fully considered**

2. Proposed Criterion



(1) Formalize Generalization Error of Representing Learning

$$h^* = \arg \min_{h \in \mathcal{H}} \sum_{i=1}^n \ell(y_i, h(x_i)) + \mathcal{R}(h)$$

General Framework of Machine Learning

The learning process is decomposed into two processes,

① Learning classifier

② Learning representation

$$h^* = g^* \circ \varphi^*$$

The final model

$$g^* = \arg \min_{g \in \mathcal{H}(\varphi^*(\mathcal{X}))} \sum_{i=1}^n \ell(y_i, g(\varphi^*(x_i))) + \mathcal{R}(g \circ \varphi^*)$$

Outer: Learning classifier

$$s.t. \varphi^* = \arg \min_{\varphi \in \Psi} P_{err}(\mathcal{H}(\varphi(\mathcal{X})))$$

Generalization Error of Representing Learning

Inner: Learning representation

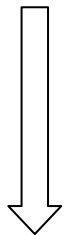
$\mathcal{H}(\varphi(\mathcal{X}))$ is the set of hyperplanes on space $\varphi(\mathcal{X})$

(2) The Upper Bound of Generalization Error



Generalization Error of Representing Learning

$$\min_{\varphi \in \Psi} P_{err}(\mathcal{H}(\varphi(\mathcal{X})))$$



- ① The upper bound of **Generalization error** is dominated by VC
- ② The **VC dimension** is dominated by the **ratio between Radius and Margin**

$$\min_{\varphi \in \Psi} \frac{R^2(\varphi(X))}{M^2(\varphi(X_+), \varphi(X_-))}$$

Not Executable

$$\text{s.t. } \begin{cases} \frac{n_{err}}{n} \leq \varepsilon \\ M(\varphi(X_+), \varphi(X_-)) = \sup_{h \in \mathcal{H}(\varphi(X))} g_m(h) \end{cases}$$

Geometric Meaning

Numerator: Radius of training set

Denominator: Margin of hyperplane

① **Theorem 2.4** (Corollary in Chapter 5.4 of (Vapnik, 1999)). With probability $1 - \eta$ one can assert that the probability that a test sample will not be separated correctly by the M -margin hyperplane has the bound $P_{err} \leq \frac{n_{err}}{n} + B_2(n, n_{err}, \eta, d_{VC})$, where $B_2(n, n_{err}, \eta, d_{VC}) = \frac{\varepsilon}{2} \left(1 + \sqrt{1 + \frac{4n_{err}}{n\varepsilon}}\right)$, $\varepsilon = 4 \frac{d_{VC}(\ln \frac{2n}{d_{VC}} + 1) - \ln \frac{\eta}{4}}{n}$, n is the number of training samples, n_{err} is the number of training samples that are not separated correctly by this M -margin hyperplane, and d_{VC} is the VC dimension in **Theorem 2.3**.

② **Theorem 2.3** (Theorem 5.1 of (Vapnik, 1999)). Let vectors $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$ belong to a sphere of radius R . Then the set of M -margin separating hyperplanes has VC dimension d_{VC} bounded by the inequality

$$d_{VC} \leq B_1(d, R, M) = \min \left(\left\lceil \frac{R^2}{M^2} \right\rceil, d \right) + 1.$$

(3) Making Margin and Radius Executable



$$\min_{\varphi \in \Psi} \frac{R^2(\varphi(X))}{M^2(\varphi(X_+), \varphi(X_-))}$$

Not Executable

s.t. $\begin{cases} \frac{n_{err}}{n} \leq \varepsilon \\ M(\varphi(X_+), \varphi(X_-)) = \sup_{h \in \mathcal{H}(\varphi(X))} g_m(h) \end{cases}$

① **Theorem 3.1.** The optimization problem (8) can be bounded by the following convex optimization problem. See Appendix A.1 for proof.

$$M^2(X_+, X_-) \geq \frac{1}{4} \min_{\alpha \in \Delta^{n_+}, \beta \in \Delta^{n_-}} \|X_+ \alpha - X_- \beta\|_2^2$$

② **Theorem 3.2.** Given a set $X \subset \mathbb{R}^d$, the squared radius of X can be computed by the following convex optimization problem. See Appendix A.2 for proof.

$$R^2(X) = \inf_{\theta \in \Delta^{|X|}} \sup_{x \in X} \|x - X\theta\|_2^2$$

- ① **Theorem 3.1 Margin** can be bounded by a convex optimization problem
- ② **Theorem 3.2 Radius** can be solved by a convex optimization problem
- ③ **Remark A.5** $M > 0 \Rightarrow n_{err}/n = 0$

$$\min_{\varphi \in \Psi} \frac{g_1(\varphi)}{g_2(\varphi)}$$

Executable

s.t. $\begin{cases} g_1(\varphi) = \min_{c \in ch(\varphi(X))} \max_{p \in \varphi(X)} \|c - p\|_2^2 \\ g_2(\varphi) = \min_{p \in ch(\varphi(X_+)), q \in ch(\varphi(X_-))} \|p - q\|_2^2 \end{cases}$

Can be Calculated Effectively and Efficiently

3. Application I: Kernel Selection



Kernel Selection Framework

Candidate set of kernel functions

- kernel type
- kernel parameter

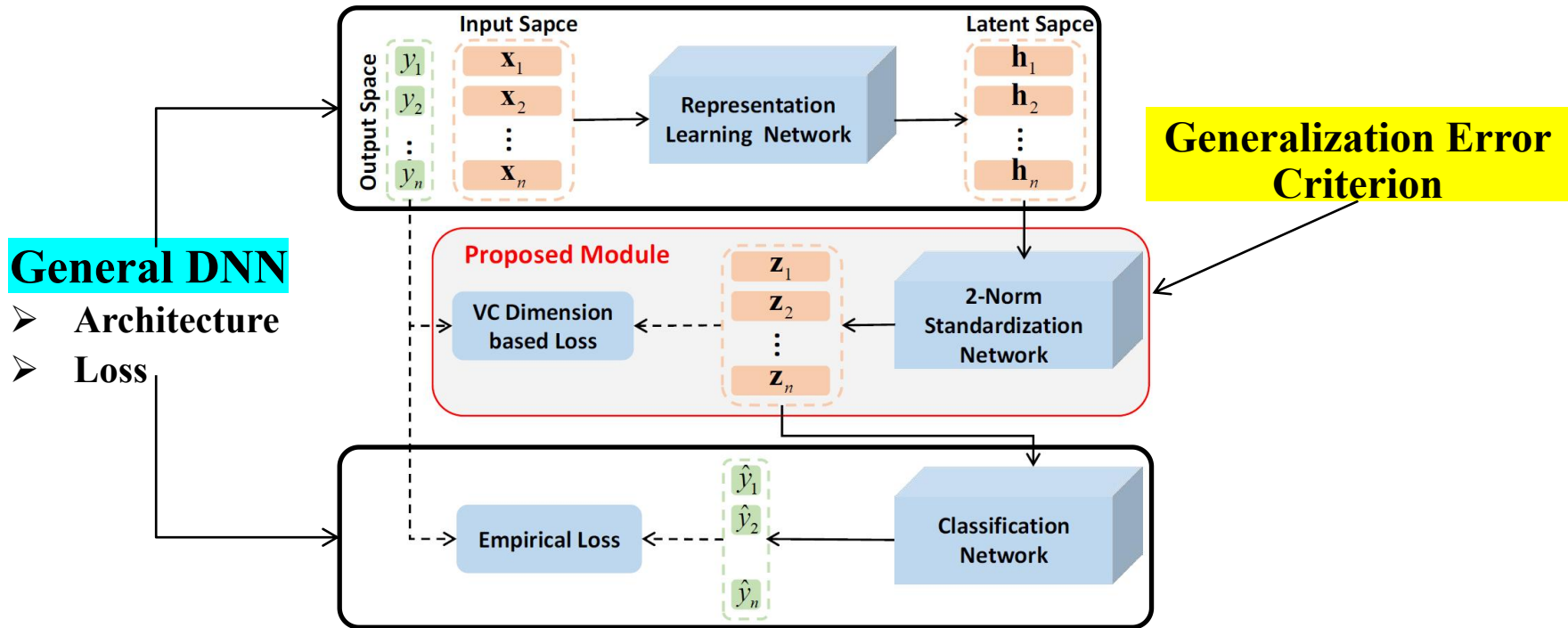
$$\min_{\varphi \in \Psi_{\text{kernel}}} f(\varphi) = 4 \frac{g_1(\varphi)}{g_2(\varphi)} \leftarrow \text{Generalization Error Criterion}$$
$$s.t. \begin{cases} g_1(\varphi) = \min_{\mathbf{u}} \max_{\mathbf{x}_i \in X} \sigma(\mathbf{u})^T \mathbf{K}^\varphi \sigma(\mathbf{u}) - 2\mathbf{K}_{[i,:]}^\varphi \sigma(\mathbf{u}) + k_{ii}^\varphi \\ g_2(\varphi) = \min_{\mathbf{v}, \mathbf{w}} \begin{pmatrix} \sigma(\mathbf{v}) \\ \sigma(\mathbf{w}) \end{pmatrix}^T \hat{\mathbf{K}}^\varphi \begin{pmatrix} \sigma(\mathbf{v}) \\ \sigma(\mathbf{w}) \end{pmatrix} \end{cases}$$

Select the kernel function with the **smallest generalization error**

4. Application II: DNN Boosting



■ DNN Boosting Framework



Using the proposed criteria to train the parameters

5. Conclusion and Outlook



1. A **crit**erion can measure the **generalization error** of **representation learning** and can be **calculated Effectively and Efficiently (Have completed)**
2. **Successful application** in kernel selection **(Have completed)**
3. **Successful application** in boosting DNN **(Have completed)**
4. A **Powerful tool** for analyzing other methods **(Be going to)**
5. **Provide Guidance** for designing new methods **(Be going to)**

Thanks!



School of Computer and Information Technology
(School of Big Data), Shanxi University
<http://cs.sxu.edu.cn/index.html>



Key Laboratory of Computational Intelligence and Chinese
Information Processing of Ministry of Education, Shanxi
University, <http://cicp.sxu.edu.cn/index.htm>

Correspondence to: Jiye Liang <lji@sxu.edu.cn>