

Singular learning theory V : symmetry and RLCT

sd15
29/4/20
①

The central statement of singular learning theory is that the Bayes generalisation error of a singular model is determined by the RLCT, and in contradistinction to the case of regular models, the RLCT is not determined by the number of parameters in the model. With a fixed class of models, the RLCT varies as the true distribution varies. In connection with this, Watanabe [W, §7.6] makes the remarkable, and somewhat cryptic, statement that

simple function \iff complicated singularities

(1.1)

complicated function \iff simple singularities.

Here "function" refers to the true distribution, and the complexity of a singularity is measured by the RLCT (complex singularities have smaller RLCTs). Since Kolmogorov complexity, and other complexity measures, lie at the heart of information theory, the above statement is potentially one of the deepest insights of statistical learning theory and information science, but it is not widely known.

In order to make Watanabe's discovery more accessible, we exhibit in this note the connection between complexity of functions and complexity of singularities in what we hope is a simple and transparent way, by emphasising the role of symmetry.

A highly symmetric function is simple, because it can be described with less information: a rotationally invariant function $f(x, y)$ may be described as a function $g(r)$. Let us explain how highly symmetric true distributions $q(y|x)$ lead to small RLCTs.

Remark Recall the RLCT is a measure of "effective number of parameters" in a model close to the most complex singularity of the set of true parameters, in the sense that in local coordinates u_1, \dots, u_d the set of true parameters is $u_1 = \dots = u_{2\lambda} = 0$ where λ is the RLCT (this is only strictly true in the "mildly singular" case, e.g. reduced rank regression). Since varying the remaining $u_{2\lambda+1}, \dots, u_d$ doesn't change the "fit" of the model, we do not count them as parameters.

We take the same setup as the "Fisher for feedforward" notes (fforw), but where $f(x, w)$ is not necessarily a ReLU network. Then

$$K(w) = \int q(y|x) q(x) \log \frac{q(y|x)}{p(y|x, w)} dx dy. \quad (1.1)$$

We do not assume the true distribution is realisable, set $W_\alpha = \{w \mid K(w) = \alpha\}$. As discussed in (sd14) beginning on p. (16), for $C \subseteq W$ compact (ignoring priors) if $W_0 \cap C \neq \emptyset$ the local RLCT $2\lambda_C$ is a measure of the effective codimension of $W_0 \cap C$ in C (very roughly). In particular, it seems reasonable to assume that the more directions at $P \in W_0 \cap C$ tangent to $W_0 \cap C$, the smaller the RLCT (this is strictly true in the "mildly singular" case from (sd14) p. (17)).

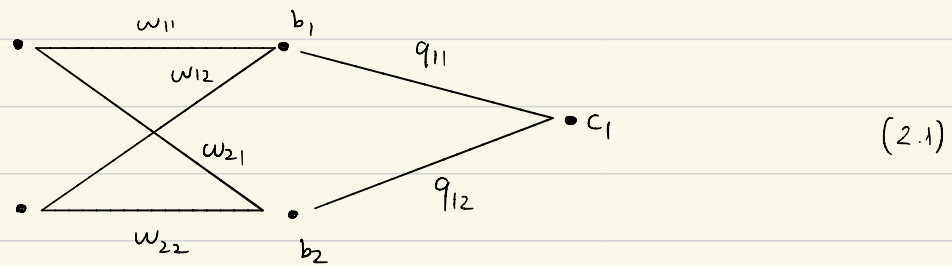
Suppose given a group G , and (not nec. continuous) actions $G \times W \rightarrow W$, written $(g, w) \mapsto g \cdot w$, and $G \times \mathbb{R}^N \rightarrow \mathbb{R}^N$, written $(g, x) \mapsto gx$.

Defⁿ We say the pairing $f: \mathbb{R}^N \times W \rightarrow \mathbb{R}^M$ is G -invariant if

$$f(gx, gw) = f(x, w) \quad (1.2)$$

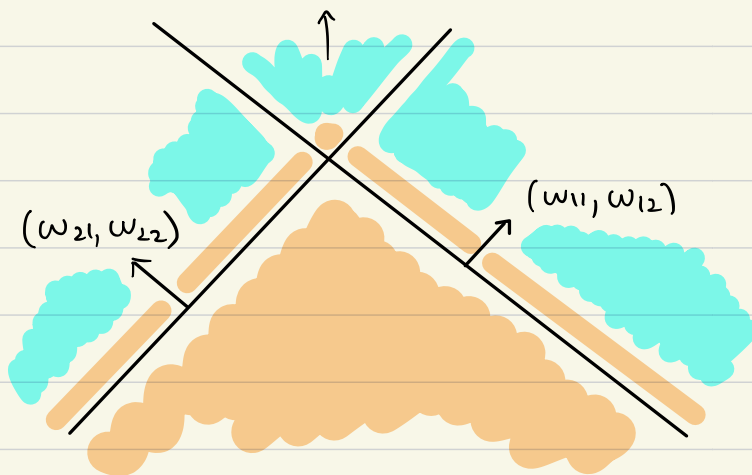
for all $x \in \mathbb{R}^N, w \in W, g \in G$ (equiv. $f(gx, w) = f(x, g^{-1}w)$ for all x, w, g).

Example 1 Consider a two layer feedforward ReLU net



with $W = \overline{B}_K(0) \subseteq \mathbb{R}^9$ for some K , $N=2$ and $M=1$, and $\omega = (\omega_{11}, \omega_{12}, \omega_{21}, \omega_{22}, b_1, b_2, q_{11}, q_{12}, c_1)$ determining for $x = (x_1, x_2)$

$$f(x, \omega) = q_{11} \text{ReLU}(\omega_{11}x_1 + \omega_{12}x_2 + b_1) + q_{12} \text{ReLU}(\omega_{21}x_1 + \omega_{22}x_2 + b_2) + c_1. \quad (2.2)$$



$$\begin{aligned} c_1 &< 0 \\ q_{11} &> 0 \\ q_{12} &> 0 \end{aligned}$$

>0

<0

Let $G = O(2)$ act on $\mathbb{R}^N = \mathbb{R}^2$ as usual, and on W by

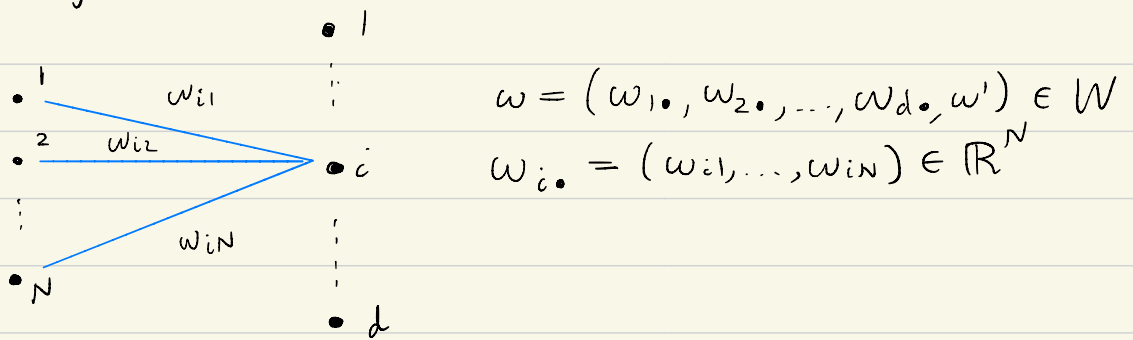
$$g \cdot \omega = (g \omega_{1\bullet}, g \omega_{2\bullet}, b_1, b_2, q_{11}, q_{12}, c_1). \quad (2.3)$$

where $\omega_{1\bullet} = (\omega_{11}, \omega_{12})^T$, $\omega_{2\bullet} = (\omega_{21}, \omega_{22})^T$. We claim that $f : \mathbb{R}^N \times W \rightarrow \mathbb{R}$ is G -invariant, with respect to these actions. Writing \langle, \rangle for the dot product this follows from

$$\begin{aligned}
 f(gx, g\omega) &= q_{11} \text{ReLU}(\langle g\omega_{1\bullet}, gx \rangle + b_1) \\
 &\quad + q_{12} \text{ReLU}(\langle g\omega_{2\bullet}, gx \rangle + b_2) + c_1 \\
 &= q_{11} \text{ReLU}(\langle \omega_{1\bullet}, x \rangle + b_1) \\
 &\quad + q_{12} \text{ReLU}(\langle \omega_{2\bullet}, x \rangle + b_2) + c_1 \\
 &= f(x, \omega).
 \end{aligned} \tag{3.1}$$

More generally:

Lemma Let $f: \mathbb{R}^N \times W \rightarrow \mathbb{R}^M$ be a ReLU network of arbitrary depth, and represent W as $W_1 \times W'$ where $W_1 = (\mathbb{R}^N)^d$ are the weights in the first layer:



We let $G = O(N)$ act on \mathbb{R}^N in the standard way, and on W by

$$\begin{aligned}
 G \times W_1 \times W' &\longrightarrow W_1 \times W' \\
 (g, (\omega_{1\bullet}, \dots, \omega_{d\bullet}), \omega') &\longmapsto (g\omega_{1\bullet}, \dots, g\omega_{d\bullet}, \omega').
 \end{aligned}$$

Then f is G -invariant.

Proof In the first layer, f computes pre-activations as

$$x \longmapsto (\langle \omega_{1\bullet}, x \rangle, \dots, \langle \omega_{d\bullet}, x \rangle). \quad \square$$

Lemma If f is G -invariant then $p(y|gx, g\omega) = p(y|x, \omega)$

Proof Clear since $p(y|x, \omega) := \frac{1}{(2\pi)^{M/2}} \exp\left(-\frac{1}{2} \|y - f(x, \omega)\|^2\right)$. \square

Proposition Let G be a group acting on \mathbb{R}^N and W such that f is G -invariant and

- (i) $q(y|gx) = q(y|x)$ for all $x \in \mathbb{R}^N, y \in \mathbb{R}^M, g \in G$
- (ii) $q(x) = q(gx)$ for all $x \in \mathbb{R}^N, g \in G$
- (iii) For $g \in G$ the action $g \cdot (-) : \mathbb{R}^N \rightarrow \mathbb{R}^N$ is smooth and $|\det(Dg)(x)| = 1$ for all $x \in \mathbb{R}^N$.

Then $K(\omega) = K(g\omega)$ for all $\omega \in W, g \in G$.

Proof

$$\begin{aligned} K(g\omega) &= \int q(y|x) q(x) \log \frac{q(y|x)}{p(y|x, g\omega)} dx dy \\ &= \int q(y|x) q(x) \log \frac{q(y|x)}{p(y|g^{-1}x, \omega)} dx dy \\ &= \int q(y|g^{-1}x) q(g^{-1}x) \log \frac{q(y|g^{-1}x)}{p(y|g^{-1}x, \omega)} dx dy \end{aligned}$$

Since $g \cdot (-) : \mathbb{R}^N \rightarrow \mathbb{R}^N$ is smooth and bijective,

$$\begin{aligned} &= \int q(y|x) q(x) \log \frac{q(y|x)}{p(y|x, \omega)} |\det(Jg)(x)| dx dy \\ &= K(\omega). \quad \square \end{aligned}$$

Corollary The level sets $W_\alpha \subseteq W$ are G -invariant.

Example 2 In the situation of Example 1, we assume $q(x)$ is an $O(2)$ -invariant distribution on \mathbb{R}^2 , e.g. a normal distribution centered at $\underline{0}$, and an example of $q(y|x)$ satisfying hypothesis (i) of the proposition is

$$q(y|x) = \frac{1}{(2\pi)^{1/2}} \exp\left(-\frac{1}{2} \|y - h(r)\|^2\right)$$

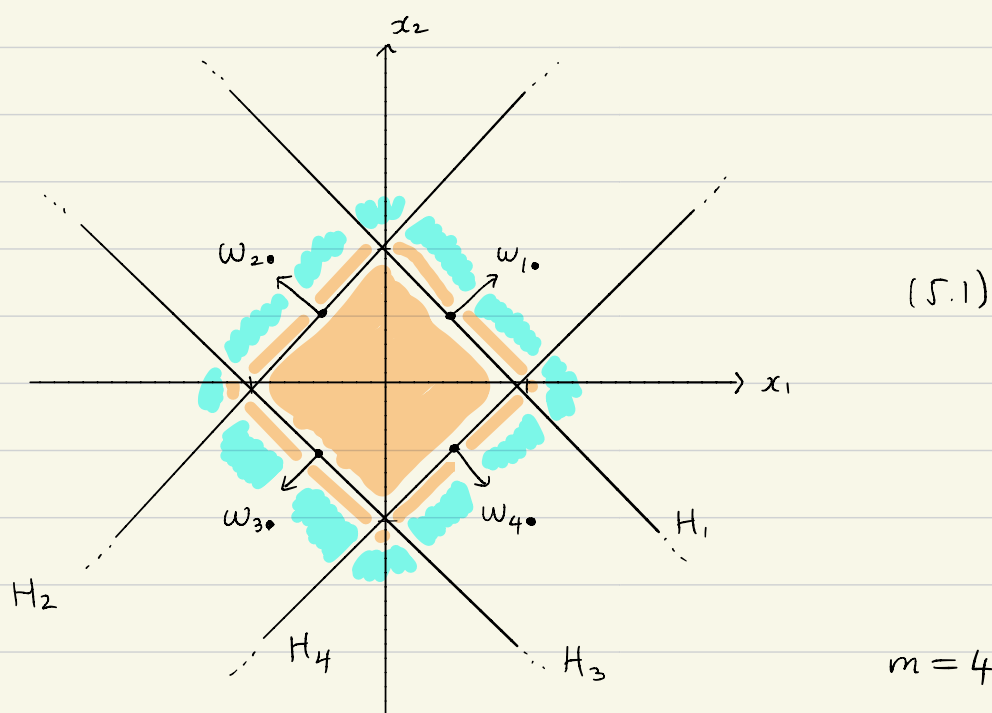
where $r = \|x\|$, for any continuous function $h: [0, \infty) \rightarrow \mathbb{R}$.

Then the proposition applies and all the level sets W_α are $O(2)$ -invariant.

The problem here is that an $O(2)$ -invariant true distribution which is nonconstant seems not to be realisable by a finite-depth ReLU network.

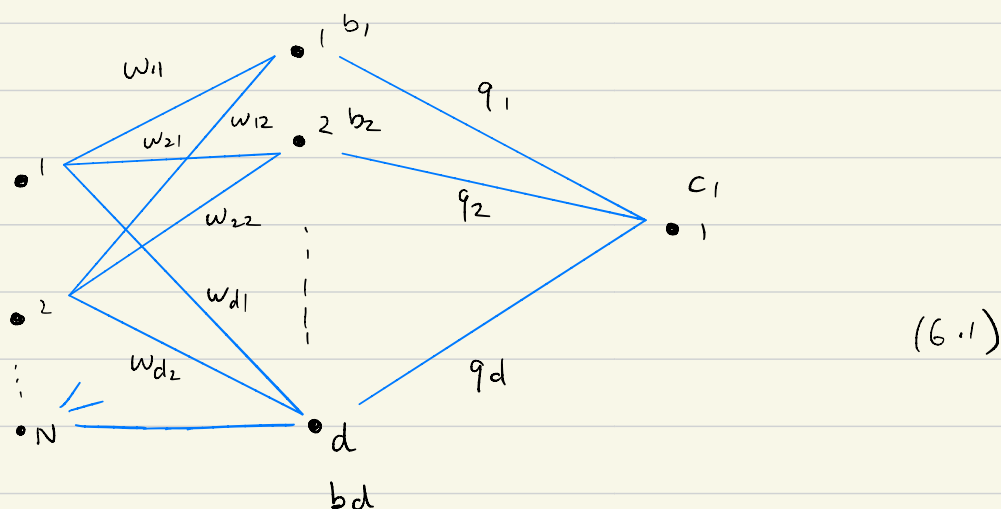
However we can consider realisable true distributions which are invariant under finite subgroups $G \subseteq O(2)$.

Example 3 Consider a depth two ReLU network $f(-, w): \mathbb{R}^2 \rightarrow \mathbb{R}$ where the hyperplanes $\langle w_i, x \rangle + b_i = 0$ (for $1 \leq i \leq d$ indexing a node in the first hidden layer) are mapped to one another by a finite subgroup $G \subseteq O(2)$ generated by rotations by $\frac{2\pi}{m}$, e.g.



This is a constraint on $w_0 \in W$ which ensures that $q(y|x) := p(y|x, w_0)$ satisfies the invariance condition of the Proposition on p. (4) for $g \in G$. So this true distribution is both G -invariant and realisable (by construction).

More generally, let $G \subseteq O(N)$ be a finite subgroup, $G = \langle g \rangle$ and suppose the network is given by weights $w_0 = ((w_{i\bullet})_{i=1}^d, b_\bullet, q_\bullet, c_1)$ as in



such that there is a permutation $\beta \in S_d$ with

$$(i) \text{ as functions } \mathbb{R}^N \rightarrow \mathbb{R} \text{ for } 1 \leq i \leq d$$

$$\langle w_{i\bullet}, gx \rangle + b_i = \langle w_{\beta(i)\bullet}, x \rangle + b_{\beta(i)} \quad (6.2)$$

$$(ii) \quad q_{\beta(i)} = q_i \text{ for } 1 \leq i \leq d.$$

Then

$$\begin{aligned} f(gx, w_0) &= c_1 + \sum_{i=1}^d q_i \text{ReLU}(\langle w_{i\bullet}, gx \rangle + b_i) \quad (6.3) \\ &= c_1 + \sum_i q_{\beta(i)} \text{ReLU}(\langle w_{\beta(i)\bullet}, x \rangle + b_{\beta(i)}) \\ &= f(x, w_0). \end{aligned}$$

Note that condition (ii) can be realised by taking all q_i equal. To reason about (i) let us assume none of the $w_{i\bullet}$ are zero vectors, so $\langle w_{i\bullet}, - \rangle : \mathbb{R}^N \rightarrow \mathbb{R}$ is surjective and $\mathbb{R}^N / K \cong \mathbb{R}$ where K is the kernel. Let $t_i \in \mathbb{R}^N$ be such that $\langle w_{i\bullet}, -t_i \rangle = b_i$. Then

$$\langle w_{i\bullet}, x \rangle + b_i = 0 \iff \langle w_{i\bullet}, x - t_i \rangle = 0$$

and (6.2) is equivalent to (writing $t_i = g g^{-1} t_i$)

$$\langle g^{-1} w_{i\bullet}, x - g^{-1} t_i \rangle = \langle w_{\delta(i)\bullet}, x - t_{\delta(i)} \rangle$$

which we can arrange by e.g. $g^{-1} w_{i\bullet} = w_{\delta(i)\bullet}$ and $g^{-1} t_i = t_{\delta(i)}$.

Lemma Let $G \subseteq O(N)$ be a finite group generated by g , and suppose $(w_{1\bullet}, \dots, w_{d\bullet}) \in (\mathbb{R}^N)^d$ and $(t_1, \dots, t_d) \in (\mathbb{R}^N)^d$ and $\delta \in S_d$ are such that

$$(a) \quad g^{-1} w_{i\bullet} = w_{\delta(i)\bullet} \quad \text{for all } 1 \leq i \leq d$$

$$(b) \quad g^{-1} t_i = t_{\delta(i)} \quad \text{for all } 1 \leq i \leq d.$$

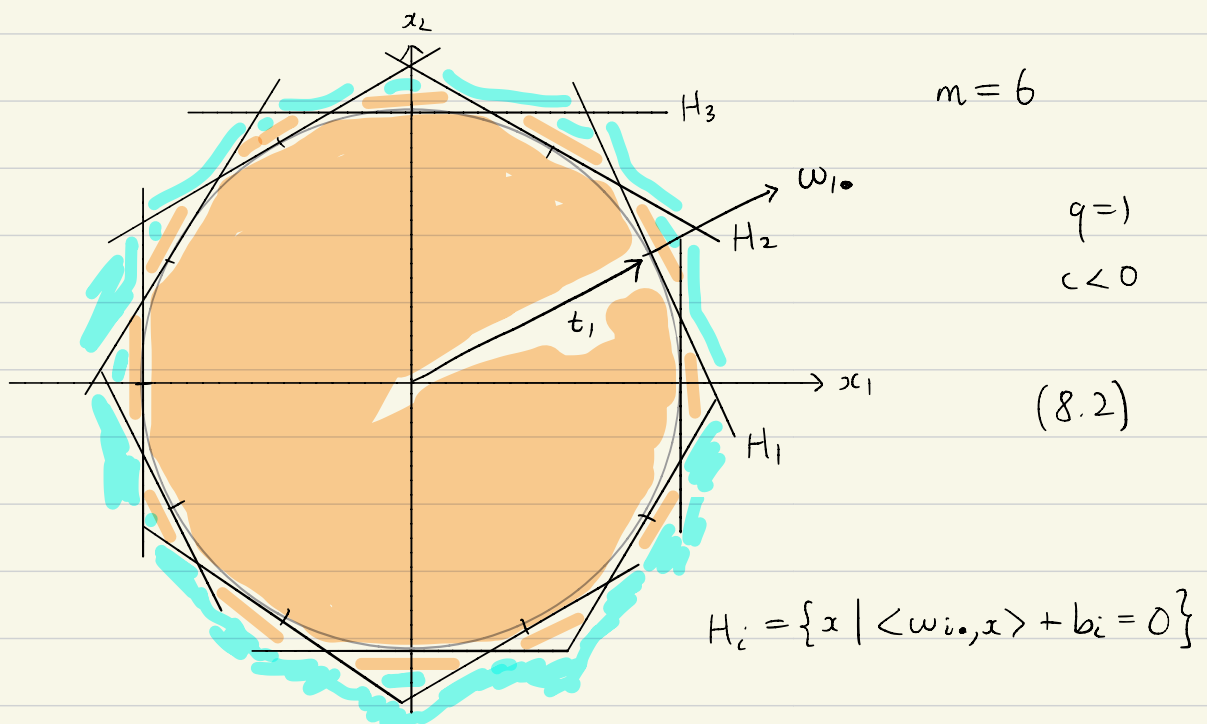
Then let $w_0 = ((w_{i\bullet})_{i=1}^d, (-\langle w_{i\bullet}, t_i \rangle)_{i=1}^d, (q)_{i=1}^d, c)$ be the parameters for a two-layer feedforward ReLU network as in (6.1), from $\mathbb{R}^N \rightarrow \mathbb{R}$, with biases $b_i = -\langle w_{i\bullet}, t_i \rangle$, and $q, c \in \mathbb{R}$ arbitrary. Then

$$f(gx, w_0) = f(x, w_0) \quad \forall x \in \mathbb{R}^N.$$

Proof Same calculation as (6.3), i.e.

$$\begin{aligned}
f(gx, w_0) &= c + q \sum_{i=1}^d \text{ReLU}(\langle w_{i\bullet}, gx \rangle + b_i) \\
&= c + q \sum_{i=1}^d \text{ReLU}(\langle w_{i\bullet}, gx - t_i \rangle) \quad (8.1) \\
&= c + q \sum_{i=1}^d \text{ReLU}(\langle g^{-1}(w_{i\bullet}), x - g^{-1}(t_i) \rangle) \\
&= c + q \sum_{i=1}^d \text{ReLU}(\langle w_{\beta(i)\bullet}, x - t_{\beta(i)} \rangle) \\
&= f(x, w_0). \quad \square
\end{aligned}$$

Example 4 To revisit Example 3 more concretely, let $g \in O(2)$ be rotation by $\frac{2\pi}{m}$ anticlockwise for some $m \geq 3$ and let $G = \langle g \rangle \cong \mathbb{Z}/m\mathbb{Z}$. Let $t_1 := g^{1/2}(1, 0)^T$ where $g^{1/2}$ is rotation by $\frac{2\pi}{2m}$, define also $w_{i\bullet} = t_i := g^i t_1$, for $1 \leq i \leq m$.



Then $g^{-1}w_{i\bullet} = g^{i-1}t_1 = w_{(i-1)\bullet}$ and $g^{-1}t_i = t_{i-1}$ where indices are read modulo m , so with β the cyclic permutation the hypotheses hold.

Taking $q(y|x) := p(y|x, w_0)$ for such w_0 constructs a \mathbb{Z}_m -invariant realisable true distribution for any $m \geq 3$. Notice that this true distribution is realisable for the ReLU architecture with d nodes in the hidden layer for any $d \geq m$ (taking some q_i 's to be zero).

Consider a two-layer ReLU network architecture $f: \mathbb{R}^2 \times W \rightarrow \mathbb{R}$ with d nodes in the hidden layer, where W is compact and $O(2)$ -invariant. For $m < d$ let $q_m(y|x)$ be the \mathbb{Z}_m -invariant realisable true distribution constructed above and λ_m its RLCT relative to some fixed prior.

Conjecture λ_m is a decreasing function of m .

Less formally

A more symmetric true distribution \Rightarrow smaller RLCT

"simpler function"

"more complicated singularity"

In their published work on RLCTs, Watanabe and collaborators tend to focus on rather simple true distributions, because it is already difficult to theoretically analyse the RLCT in these cases. However the deep idea on p. ① is best illustrated with more interesting true distributions, and experimental approximation to the RLCT.

Remark We justify the application of singular learning theory to ReLU networks by the "soft ReLU trick".

Remark The idea of studying figures in the plane in connection with properties of neural networks is inspired by Minsky & Papert's book "Perceptrons: an introduction to computational geometry".