This note is a continuation of "Singular Learning Theory 5" Solis, where a particular class of models based on feedforward ReLU networks were defined. In this note we theoretically derive the RLCT of these models paired with the symmetric true distributions.

Recall from (6.1) of Gals that we consider a two-layer ReLU network



and the corresponding function

$$f : \mathbb{R}^2 \times \mathbb{W} \longrightarrow \mathbb{R} \tag{1.2}$$

$$f(x,w) = c_1 + \sum_{i=1}^{d} q_i \operatorname{ReLU}(\langle w_{i}, x \rangle + b_i)$$

where 
$$w = (w_{10}, ..., w_{d0}, b_{1}, ..., b_{d}, q_{1}, ..., q_{d}, C_{1})$$
, and  $W \subseteq \mathbb{R}^{41}$   
is compact. We denote by Hi the following subspace

$$H_{i} = \left\{ x \in \mathbb{R}^{2} \mid \langle w_{i}, x \rangle + b_{i} = 0 \right\}$$
(1.3)

We assume the distribution q(y|x) is realisable, and that a prior  $\mathcal{P}$  is fixed which is nonzero on the set of true parameters.

The model class is p(y|x, w) where (see "Fisher for feedforward")

$$p(y|x,\omega) := \frac{1}{(2\pi)^{1/2}} \exp\left(-\frac{1}{2} \|y - f(x,\omega)\|^2\right)$$
(2.1)

and the Kullback - Leibler distance to the true distribution is

$$K(\omega) = \int q(y|x) q(x) \log \frac{q(y|x)}{p(y|x,\omega)} dxdy \qquad (2.2)$$

where q(x) is a chosen probability density function on  $\mathbb{R}^2$ . We are first of all interested in the set of true parameter

$$W_{o} = \left\{ w \in W \mid K(w) = 0 \right\}$$
(2.5)

which is nonempty by hypothesis. The twe distribution we consider here is inspired by, but different from, that in SdIS. The reason for the change is that the  $\mathbb{Z}/m\mathbb{Z}$  symmetry no longer appears central.

<u>Hypothesis I</u> we suppose given an integer  $I \le m \le d$  and lines  $H_{i}^{(0)}$ ...,  $H_{m}^{(0)} \le \mathbb{R}^{2}$ together with chosen  $W_{i}^{(0)} \in \mathbb{R}^{2}$  and  $b_{i}^{(0)} \in \mathbb{R}$  for  $I \le i \le d$  such that  $H_{i}^{(0)}$ is the line determined by  $W_{i}^{(0)}$ ,  $b_{i}^{(0)}$  as in (1.3) (note that since  $H_{i}$  is a line, this means  $W_{i}^{(0)} \neq 0$  for  $I \le i \le d$ ) and  $H_{i}^{(0)} \neq H_{j}^{(0)}$  whenever  $i \ne j$ . Then

$$q(y|x) := p(y|x, \omega^{(0)}) \qquad (2.4)$$

where  $\omega^{(0)}$  is the weight vector (assumed to be in the interior of W)

$$\omega^{(o)} = \left( \begin{array}{ccc} \omega_{1,0}^{(o)}, \dots, \omega_{m,0}^{(o)}, 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} \omega_{1,0}^{(o)}, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0, \dots, 0, \\ d-m \end{array} \right) \underbrace{d-m}^{(o)} \left( \begin{array}{ccc} 0$$

5418) 8

Explicitly

$$q(y|x) = \frac{1}{(2\pi)^{1/2}} \exp\left(-\frac{1}{2} \|y - \sum_{i=1}^{m} \operatorname{ReLU}(\langle w_{i\bullet}^{(o)}, x \rangle + b_{i}^{(o)}) \|^{2}\right) \quad (3.1)$$

so the two distribution is a generic model determined by a ReLU network of depth two with m hidden nodes (the values of c and the qc not being important).

Definition Given w, w' EW we define

$$K(\omega',\omega) = \int p(y|x,\omega') q(x) \log \frac{p(y|x,\omega')}{p(y|x,\omega)} dx dy$$

so that  $K(\omega) = K(\omega^{(\circ)}, \omega)$ .

We calculate

Fix

$$K(w',w) = \int q(x) \frac{1}{(2\pi)^{1/2}} \exp\left(-\frac{1}{2} || y - f(x,w') ||^{2}\right)$$

$$\cdot \log \frac{\exp\left(-\frac{1}{2} || y - f(x,w') ||^{2}\right)}{\exp\left(-\frac{1}{2} || y - f(x,w) ||^{2}\right)} dxdy$$

$$= \frac{1}{(2\pi)^{1/2}} \int q(x) \exp\left(-\frac{1}{2} || y - f(x,w) ||^{2}\right)$$

$$\cdot \frac{1}{2} \left\{ || y - f(x,w) ||^{2} - || y - f(x,w') ||^{2} \right\} dxdy$$

$$x \in \mathbb{R}^{2} \text{ and let } u = y - f(x,w'), \quad a = f(x,w) - f(x,w') \quad s_{0} \text{ that}$$

$$K(w',w) = \frac{1}{2(2\pi)^{1/2}} \int q(x) K(w',w,x) dx \qquad y^{2} - 2au + a^{2} - \alpha^{2}$$

$$K(w',w,x) = \int \exp\left(-\frac{1}{2}u^{2}\right) \cdot \left\{ (u - a)^{2} - u^{2} \right\} du$$

$$= a^{2} \int \exp\left(-\frac{1}{2}u^{2}\right) du - 2a \int u \exp\left(-\frac{1}{2}u^{2}\right) du$$

Hence by standard techniques (see "Fisher for feedforward")

$$K[\omega',\omega,\varkappa) = \int 2\pi q^2 \qquad (4.1)$$

Hence

$$K(\omega', \omega) = \frac{1}{2} \int q(x) \| f(x, \omega) - f(x, \omega') \|^2 dx \qquad (4.2)$$

Note that since f(x,w) is piecewise-linear, on partitions determined by at most a hyperplanes, to determine if f(x,w) = f(x,w') on all of  $\mathbb{R}^2$  it suffices to compare their values in some compact neighbourhood of O. Without attempting to make this precise (at the moment) we assume q(x) is positive in a big enough neighbourhood of O so that this comparison suffices for all  $w \in W$ .

Hypothesis I for all 
$$w \in W$$
,  $q(x)$  is such that  $K(w) = O$  if and only if

$$f(x, w) = f(x, w^{(0)})$$
 for all  $x \in \mathbb{R}^2$ .

So (4.2) should be read as an assertion that there are many natural choices for q(a) satisfying Hypothesis II. Hence

$$W_{o} = \left\{ w \in W \mid f(x,w) = f(x,w^{(o)}) \text{ for all } x \in \mathbb{R}^{2} \right\}$$
(4.3)

which reduces the problem of classifying true parameters to the problem of classifying weights w & W up to functional equivalence of the corresponding ReLU network. For this we use [PL] as a starting point.

<u>Remark</u> What we have said so far generalises to ReLU networks of arbitrary depth computing functions  $\mathbb{R}^{N} \longrightarrow \mathbb{R}^{M}$ .



Fold sets following the notation of [PL] we define the fold set of a continuous piecewise-linear function  $f: \mathbb{R}^2 \longrightarrow \mathbb{R}$  by

 $\mathcal{F}(f) = \{x \in \mathbb{R}^2 \mid f \text{ is not clifferentiable at } x\}$ 

<u>Def</u><sup>n</sup> Given a piecewise-linear  $f: \mathbb{R}^2 \longrightarrow \mathbb{R}$  let  $\{U_{\mathcal{A}}\}_{\mathcal{A}\in\mathcal{A}}$  be the domains in  $\mathbb{R}^2$ such that  $U_{\mathcal{A}} \cup_{\mathcal{A}} = \mathbb{R}^2 \setminus F(f)$ . Then for each  $\alpha$  there are unique  $\omega^{\alpha} \in \mathbb{R}^2$ and  $b^{\alpha} \in \mathbb{R}$  such that  $f(x) = \langle \omega^{\alpha}, x \rangle + b^{\alpha}$  for all  $d \in \mathcal{O} \alpha$ .

Example By hypothesis  $F(f(x, w^{(\circ)})) = \bigcup_{i=1}^{m} H_{i}^{(\circ)}$ 

Note that in general, for  $w \in W$  some of the Hi may be empty, or all of  $|\mathbb{R}^2$ , and even if they are all lines the foldset may be a proper superset of  $U_i$  Hi, as the following example shows.

Example Consider the network  $f(x, w) = c + q \operatorname{ReLU}(wx_1 + w'x_2 + b)$ 



whose foldset is, assuming  $q \neq 0$ , the line  $w \propto_1 + w' \propto_2 + b = 0$ .

Next we consider the simplest case, where m (the number of hyperplanes determining the twe distribution q(y|x)) is equal to cl (the number of hyperplanes "available" in the model).

Case 
$$m=d$$
 Suppose  $f(x,w) = f(x,w^{(n)})$  as functions on  $\mathbb{R}^3$ . Then

$$\mathcal{F}(f(x,\omega)) = H_1^{(0)} \cup \cdots \cup H_m^{(0)}$$

which implies that (writing  $\omega$ , 9, b, c for the components of  $\omega$ ) none of the  $q_i$  are zero, and none of  $\omega_i$ . are zero for  $1 \le i \le m$ . Further there is a permutation  $3 \le S_3$ such that  $H_i = H_{si}^{(0)}$  for  $1 \le i \le m$ .

Lemma For  $w, w' \in \mathbb{R}^2$  and  $b, b' \in \mathbb{R}$  we set

$$H = \left\{ x \in \mathbb{R}^2 \mid \langle w, x \rangle + b = 0 \right\}$$
$$H' = \left\{ x \in \mathbb{R}^2 \mid \langle w', x \rangle + b' = 0 \right\}.$$

Suppose H, H' are both lines. Then H = H' if and only if there exists  $\lambda \in \mathbb{R} \setminus \{0\}$  such that  $\omega = \lambda \omega'$ ,  $b = \lambda b'$ .

Proof One direction is clear. Since  $w \neq 0$  we can find unique  $t \, s.t. \langle w, t \rangle = -b$ , and similarly  $\langle w', t' \rangle = -b'$ , with t a multiple of wand t' a multiple of w'. Then  $t \in H$ ,  $t' \in H'$  and  $\langle w, x - t \rangle = 0 \iff x \in H$  $\Leftrightarrow x \in H'$  $\Leftrightarrow \langle w', x - t' \rangle = 0$ 

Hence  $\langle w', t - t' \rangle = 0$  and  $\langle w, t - t' \rangle = 0$ , so if w, w' are linearly independent t = t', which is a contradiction since  $t \in span(w)$ ,  $t' \in span(w')$ . Hence  $\{w, w'\}$  is linearly dependent, say  $w = \lambda w'$ . Then  $-b = \langle w, t \rangle = \langle w, t' \rangle = \lambda \langle w', t' \rangle = -\lambda b'$  so  $b = \lambda b'$ .



It follows that for  $1 \le i \le m$  there exists a (unique)  $\lambda_i \in \mathbb{R}_{>0}$  and  $\in_i \in \mathbb{Z}_2$  s.t.

$$\begin{aligned}
\omega_{i \bullet} &= \lambda_{i} (-1)^{\epsilon_{i}} \omega_{\epsilon_{i}}^{(\bullet)} \\
b_{i} &= \lambda_{i} (-1)^{\epsilon_{i}} b_{\epsilon_{i}}^{(\bullet)}
\end{aligned} \tag{7.1}$$

Hence

$$f(x,\omega) = c_{1} + \sum_{j=1}^{m} q_{i} \operatorname{ReLU}(\langle \omega_{i}, x \rangle + b_{i})$$

$$= c_{1} + \sum_{i=1}^{m} q_{i} \operatorname{ReLU}(\lambda_{i}(-1) \{\langle \omega_{6i}, x \rangle + b_{3i}^{(\circ)} \}) \quad (7-2)$$

$$= c_{1} + \sum_{i=1}^{m} q_{i} \lambda_{i} \operatorname{ReLU}((-1)^{e_{i}} \{\langle \omega_{6i}, x \rangle + b_{3i}^{(\circ)} \})$$

Using  $\omega$  we may assign an orientation to the normal bundle of the 1-strata of  $\mathcal{F}(f(z,\omega))$  in  $\mathbb{R}^2$ , according to the direction of the velevant  $\omega_{i\bullet}$ . More precisely, let  $x \in S \subseteq \mathcal{F}(f(z,\omega))$  be a point in a 1-strata S, which is contained in  $H_{i\circ}$ . Then  $\omega_{i\bullet}$  gives the direction in which the corresponding ReLU is active ".



Picking  $x \in U = U_+ \cup U_-$  small enough we have  $f(x, \omega) = \langle \omega^{\pm}, x \rangle + b^{\pm}$  for  $x \in U^{\pm}$  where  $\omega^{\pm} - \omega^{-} = q_i \, \omega_i \, .$  The point here is that  $\omega^{\pm} - \omega^{-}$  may be inferred from the function  $f(x, \omega)$  by looking at its values near x, and hence so may  $q_i \omega_i$ .



The line R is divided into regions  $U_{\alpha} = (-\infty, -1)$ ,  $U_{\beta} = ($ 

This means that for each 15 is m we are in one of two cases:

(a) 
$$\epsilon_i = 0$$
, so  $\omega_i$ ,  $\omega_{\epsilon_i}$  are the same clirection, and  $q_i \omega_i = q_i^{(\circ)} \omega_{\epsilon_i}$ , so

$$\omega_{i\bullet} = \frac{q_i^{(o)}}{q_i} \omega_{z_i,\bullet}^{(o)} \tag{9.1}$$

(b)  $\in_i = 1$  so  $w_{i_0}, w_{i_0}$  are in opposite directions and  $q_i w_{i_0} = -q_i^{(\circ)} w_{\delta_i}$ , so

$$\omega_{i\bullet} = -\frac{q_i^{(0)}}{q_i} \omega_{\delta_{i/\bullet}}^{(0)} \tag{9.2}$$

Together this shows that for  $| \le \hat{i} \le m$  (noting  $q_i^{(6)} = 1$ )

$$\omega_{i} = (-1)^{\epsilon_{i}} \frac{q_{i}^{(0)}}{q_{i}} \omega_{s_{i}}^{(0)}$$
(9.3)

Companison with (7.1) shows that 
$$q_i$$
 must be positive, that  $\lambda_i = /q_i$  and

$$b_{i} = (-1)^{\epsilon_{i}} \frac{9i^{(0)}}{9i} b_{2i}^{(0)}$$
(9.4)

As the example on the previous page shows, it is not necessarily the case that all E: are zero. But these may be nonzero only under special conditions:

Lemma Set 
$$F = \{i \mid \epsilon_i = 1\}$$
. Then  $\sum_{i \in F} \omega_{\delta_i} = 0$ .

<u>Proof</u> With the notation of p. (1) let  $U \propto be a \text{ domain and let } f_i^{\alpha} be + 1 \text{ if}$ this domain is in the active region for the ReLU with vector  $w_i^{(0)}$  in the network  $w^{(0)}$  computing  $f(x, w^{(0)})$  (so  $\langle w_{i,0}^{(0)}, x \rangle + b_i^{(0)} > 0$ ), and zero otherwise. Let  $\overline{J_i}^{\alpha}$  be +1 if  $\overline{J_i}^{\alpha} = 0$  and 0 otherwise. Then

$$\sum_{i} \delta_{ii}^{\alpha} w_{ii}^{(0)} = \sum_{i \notin F} \delta_{ii}^{\alpha} q_{i}^{i} w_{i,i} + \sum_{i \in F} \delta_{ii}^{\alpha} q_{i}^{i} w_{i,i}. \qquad (9.5)$$

Hence  

$$\begin{aligned}
\sum_{i \notin F} \delta_{\delta i}^{A} \omega_{\delta i}^{(o)} + \sum_{i \in F} \delta_{\delta i}^{A} \omega_{\delta i}^{(o)} \\
&= \sum_{i \notin F} \delta_{\delta i}^{A} \omega_{\delta i}^{(o)} + \sum_{i \in F} \delta_{\delta i}^{A} (-\omega_{\delta i}^{(o)}) \\
&= \sum_{i \notin F} \delta_{\delta i}^{A} \omega_{\delta i}^{(o)} + \sum_{i \in F} \delta_{\delta i}^{A} (-\omega_{\delta i}^{(o)}) \\
&= \sum_{i \in F} \left[ \delta_{\delta i}^{A} + \overline{\delta}_{\delta i}^{A} \right] \omega_{\delta i}^{(o)} = 0 \text{ as claimed} \cdot \prod
\end{aligned}$$

$$\frac{E_{cample}}{\omega_{i}} \text{ (onsider } f(x_{i}\omega) \text{ determing hyperplanes as shown, } with vectors \\
&= \omega_{i}, \omega_{2,i}, \omega_{3}. \text{ In all six regions some ReLU is active} \\
&= \sum_{i \notin F} \left( \sum_{j \in F} (-\omega_{j}) - \sum_{j \in$$

Suppose the function f(x, w) is known, together with the  $q_{\hat{x}}, w_{\hat{i}}$  and  $b_{\hat{c}}$ . Then  $c_1$  may be found by evaluating

$$c_{1} = f(x, w) - \sum_{i=1}^{m} q_{i} \operatorname{ReLU}(\langle w_{i}, x \rangle + b_{i})$$

at any point of  $\mathbb{R}^2$ . One way of putting this is that if you know the decision boundaries and orientations, you know how qw and qb vary across boundaries and hence for any  $x \in \mathbb{R}^2$  you can determine  $\sum_j q_j \ b_j + c_1$  (jranging over active indices) by examining f(z, w) (i.e. this is  $b^\infty$  in the notation of p.  $\mathcal{P}$ ) and then subtract  $\sum_j q_j \ b_j$ .

$$\begin{pmatrix} \left(\lambda_{i}, q_{i}\right)_{i=1}^{m}, \delta, \epsilon \end{pmatrix} \longmapsto \begin{pmatrix} \left(\lambda_{i}(-1)^{\epsilon_{2i}} \omega_{si_{j}}^{(o)}\right)_{i=1}^{m}, & (||.|) \\ \left(\lambda_{i}(-1)^{\epsilon_{2i}} b_{bi}^{(o)}\right)_{i=1}^{m}, & (||.|) \\ \sum_{\epsilon_{i}=1}^{\infty} b_{i}^{(o)} \end{pmatrix}$$



Proof The function computed is (writing  $\epsilon'_{c} = \epsilon_{bc}$ )

$$\sum_{\epsilon_{i}=1}^{6} b_{i}^{(\circ)} + \sum_{i=1}^{m} \operatorname{ReLU}\left(\left(-1\right)^{\epsilon_{i}}\left\{\left\langle \omega_{\delta_{i}}^{(\circ)}, x \right\rangle + b_{\delta_{i}}^{(\circ)}\right\}\right)$$
(12.1)

This has the same foldset as  $f(x, w^{(0)})$  and so it suffices to show that for each open domain  $U_{\alpha}$  both (11.2) and  $f(x, w^{(0)})$  determine the same  $w^{\alpha}$ ,  $b^{\alpha}$ . Let  $\delta_i^{\alpha}$ ,  $\overline{\delta}_i^{\alpha}$  be as on p. (a) (so they are referring to  $f(x, w^{(0)})$ ). Then the "slope"  $w^{\alpha}$  computed by (11.2) is  $(F = \{i \mid e_i' = 1\})$ 

$$\sum_{i \notin F} \delta_{\lambda i}^{\alpha} \omega_{\lambda i}^{(0)} + \sum_{i \in F} \overline{\delta}_{\lambda i}^{\alpha} \left(-\omega_{\lambda i}^{(0)}\right) \qquad (12.2)$$

whereas the slope computed by  $f(x, w^{(0)})$  is  $\sum_{i} \delta_{zi}^{*} w_{zi}$ . However by hypothesis  $\sum_{i \in F} \left[ \delta_{zi}^{*} + \overline{\delta_{6i}^{*}} \right] w_{si}$ ,  $= \sum_{i \in F} w_{bi}^{(0)} = 0$  so these are equal. The "intercept" be computed by (11.2) is

$$\sum_{i \in F} b_{2i}^{(o)} + \sum_{i \notin F} \delta_{2i}^{d} b_{2i}^{(o)} + \sum_{i \in F} \overline{\delta}_{2i}^{d} (-b_{6i}^{(o)})$$

$$= \sum_{i \in F} \delta_{3i}^{d} b_{3i}^{(o)} + \sum_{i \notin F} \overline{\delta}_{2i}^{d} b_{6i}^{(o)} \quad (12.3)$$

$$= \sum_{i \in F} \delta_{3i}^{d} b_{3i}^{(o)}$$

which is the intercept computed by  $f(x, w^{(n)})$ . Hence  $f(x, w^{(n)}) = f(x, w)$ , where w is the vector on the RHS of (11.1). Twjectivity is clear, since the hyperplanes giving the decision boundaries are distinct so 3 may be recovered, and thus also  $\in$ . Sujectivity follows from what we have said above.  $\Box$  <u>Definition</u> For  $3 \in S_m$ ,  $\epsilon \in P$  let  $W_0^{3,\epsilon} \subseteq W_0$  denote the image of  $\mathcal{Y}$  restricted to  $X^m \times \{2\} \times \{\epsilon\}$ .

The subset  $W_0^{6, \epsilon} \subseteq W_0$  is a submanifold and  $W_0$  is the disjoint union over all 3,  $\in$  of the  $W_0^{6, \epsilon}$ . Let  $C^{3, \epsilon}$  be a compact neighbourhood of  $W_0^{6, \epsilon}$  in  $W_0$ which does not meet any other component  $W_0^{6', \epsilon'}$ . Recall the local RLCT of p. (G) (d)4.

<u>Definition</u> Let  $\lambda^{6,\epsilon}$  denote the learning coefficient of (P, 9, f) restricted to  $C^{b,\epsilon}$ .

Recall that W has dimension 3m+1.

Proposition  $\lambda^{\delta, \epsilon} = m + \frac{1}{2}$ , for all  $\delta, \epsilon$ .

<u>Boof</u> Since  $W_{o}^{\delta,\epsilon}$  is a submanifold of dimension m [W, Remark 7.3] gives the inequality  $\lambda^{6,\epsilon} \leq m + \frac{1}{2}$  directly. But for equality we need to say more. Set  $\Omega_{-} = C^{\delta,\epsilon}$  and let  $A_{\mathcal{R}}$  be the sheaf of real analytic functions on  $\mathcal{N}$  and  $\mathcal{I}$  the ideal sheaf generated by K so that the shucture sheaf of  $W_{o}^{\delta,\epsilon}$  is  $A_{\mathcal{R}}/\mathcal{I}$ . We claim that  $W_{o}^{\delta,\epsilon}$  is a regular real-analytic variety. This follows from the  $(R_{>o})^{m}$ -action on  $\mathcal{N}$  being transitive on  $W_{o}^{\delta,\epsilon}$  since the singular locus is closed (in whatever topology you like). Hence at all points of  $W_{o}^{\delta,\epsilon}$  the RLCT is just half the codimension [L, Theorem 5.1].  $\Box$ 

Note that the structure sheaf on  $W_0^{\delta_1 \epsilon}$  as a real-analytic variety is not determined by  $W_0^{\delta_1 \epsilon}$  as a <u>subject</u> and it is the structure sheaf which determines which points of this subvariety are regular (at a regular point the RLCT is the codimension). Since we have not analysed K we cannot say much about  $\mathcal{O}_{W_0^{\delta_1 \epsilon}}$  directly. However the singular locus is closed [H, Theorem 5.3] so a generic point of  $W_0^{\delta_1 \epsilon}$  must be regular. Since  $(R_{>0})^m$  acts transitively no point is special, so they must all be generic hence regular.



## Corollary The global RLCT of (P,9,9) is m+2.

## References

[W] S. Watanabe "Algebraic geometry and statistical learning"

[PL] M. Phuong, C. H. Lampert "Functional vs. parametric equivalence of ReLU networks" ICLR 2020.

[L] S. Lin "Vseful facts about RLCT" 2012.

[H] R. Hartshorne ('Algebraic geometry '