Singular learning theory IV : local RLCT

Throughout [W] means Watanabe's book. We work in the setting of [W, §6.1] and assume Fundamental Conditions (I), (II) with s = 2. (a somewhat informal, but more readable, account is given in [W, §1.4]). In particular q (x) is realisable, and W is as in [W, Defn 6.3]. We aim to clarify (a) the role of local RLCTs (to be defined) and (b) the role of the prior in RLCTs.

Setup Recall the partition of parameter space of [w, Theorem 6.5]. For some E>O we replace W by

$$W_{\epsilon} = \left\{ w \in W \mid K(w) \leq \epsilon \right\}$$

which is legitimate since we are only concerned with the RLCT, in this note. Then $\mathcal{M} := g^{-1}(W_{\mathcal{E}})$ is covered by a finite set

$$\mathcal{M} = \mathcal{U}_{\alpha} \mathcal{M}_{\alpha}$$

where the Ma are a very particular kind of set (not open) given in Remark 2.14, and constructed as follows: for each $p \in \mathcal{M}$ choose an open subject of a coordinate chart centered on p of the form

$$O_{p}(b) := (-b,b)^{a} = \{ u = (u_{1},...,u_{d}) \mid (u_{i}| < b | \le i \le d \}$$
 ((.))

where the dimension of W (hence also \mathcal{M}) is d. The construction of Remark 2.14 covers \mathcal{M} by $\{O_p(b)\}_{p\in\mathcal{M}}$ uses compactness to choose a finite subcover. For each p chosen, it then writes (using again local wordinates) 13/4/20 (scl14)

()

$$O_{p}(b) = [(-b,0] \cup [0,b)]^{d}$$

= (-b,0] × (-b,0] × ... × (-b,0]
 $\cup (-b,0] \times ... - - \times [0,b] \cup ...$
= $\bigcup_{i=1}^{2^{d}} M_{i}$ where $M_{i} \cong [0,b]^{d}$.

As p ranges over the finite list of chosen points, and i over indices for these
sequences of "posifive/negative" choices for local coordinates up around p,
we enumerate some finite ret of subsets
$$Ma \subseteq M$$
 such that $M = Ua Ma$.
There is an associated family of functions $\mathcal{G}_{\alpha}(u)$ with supp $\mathcal{G}_{\alpha} \subseteq Ma$
 $(\mathcal{G}_{\alpha}(u))$ being C^{ω} by the contal partition of unity tricks using the
prescribed $\mathcal{G}_{\alpha}^{(b)}$ rather than $\equiv 1$ check) such that for any integrable H,

$$\int_{W} H(w) dw = \int_{M} H(g(u)) |g'(u)| du \qquad (2.1)$$
$$= \sum_{\alpha} \int_{M_{\alpha}} H(g(u)) \partial_{\alpha} (u) |g'(u)| du$$

In [W, Theorem 6.5]; + continues to say that in each $M_{\mathcal{X}}$ (keeping in mind $0 \le u_i \le b$) there is a C^{∞} function $\phi(u)$ such that

$$\kappa(g(u)) = u^{2k} = u^{2k_1} \cdots u^{2k_d} \qquad k_i \in \mathbb{N}$$

$$\mathcal{G}(g(u)) |g'(u)| = \phi(u)u^{h} = \phi(u)u^{h} \dots u^{h} d \qquad h \in \mathbb{N}$$
(2.2)

Hence (2.1) may be further refined to

$$\int_{W} H(w) \mathcal{G}(w) dw = \sum_{\alpha} \int_{M_{\alpha}} H(\mathcal{G}(u)) \phi_{\alpha}^{*}(u) u^{h} du$$
where $\phi_{\alpha}^{*}(u) = \mathcal{E}_{\alpha}(u) \phi(u)$. (2.3)

Calculating RLCTs consider the zeta function of the triple (P, 2, Y, W)

$$\zeta(z) = \int_{W}^{r} K(w)^{2} f(w) dw \qquad (Re(z) > 0) \qquad (3.1)$$

2 cl 12

which can be analytically continued [W, Thm 6.6] to a unique meromorphic function on the entire complex plane whore poles are all real, negative and rational. The learning coefficient λ is by clef^N such that $-\lambda$ is the largest pole of $\mathcal{Z}(z)$. Here $\kappa(w)$ is defined by

$$K(\omega) = \int q(x) \log \frac{q(x)}{p(x \mid \omega)} dx \qquad (3.2)$$

Note that the learning wefficient depends, in principle, on the prior \mathcal{I} , whereas the <u>Real Log (anonical Threhold</u> (RLCT), defined in [W, Def^{*}2.7] for the pair (W, K) is associated to $\{ \omega \mid K(\omega) = 0 \}$ independently of the prior. As we will see, and as stated in [W, Remark 6.7(1)] provided that the prior is positive on $\{ \omega \mid K(\omega) = 0 \}$ (actually a weaker condition suffices) the learning wefficient is equal to the RLCT, and in particular is independent of the prior.

The first remark (proof of [W, Theorem 6.6]) is that for Re(2) > 0

and the second summand extends to a holomorphic function on the entire complex plane, so in analysing λ we may restrict to integrals over $W \in a$ s above.

Then using (2.3)

$$\overline{\zeta}_{1}(z) = \sum_{\alpha} \int_{M_{\alpha}} u^{2k z} u^{h} \phi^{*}_{\alpha}(u) du \qquad (4.1)$$

Note that b(u) in the statement of [W, Thm 2.3] is real analytic, so $\phi(u)$ may be taken real analytic provided S is (and we assume so). In this case, by shrinking the Ma if necessary, we may awange $\phi_a^*(u)$ to actually have an expansion as

$$\phi_{\alpha}^{*}(u) = \sum_{\substack{i \leq N \\ i \leq N}} a_{j}^{(\alpha)} u^{j} + R_{N}(u). \qquad (4.2)$$

Therefore in the region Re(z) > 0, taking Na sufficiently large,

$$\hat{\boldsymbol{\zeta}}_{1}(\boldsymbol{z}) = \sum_{\boldsymbol{\alpha}} \sum_{\substack{j \mid \leq N_{\boldsymbol{\alpha}}}} \int_{M\boldsymbol{\alpha}} a_{j}^{(\boldsymbol{\alpha})} u^{2\boldsymbol{k}\boldsymbol{z}} u^{\boldsymbol{k}} u^{\boldsymbol{j}} d\boldsymbol{u} \quad (4.3)$$

+ terms with holomouphic extension to all of (.

Since
$$M_{\alpha} = [0, b]^{d}$$
 this is

$$= \sum_{\alpha} \sum_{\substack{j \mid \leq N \\ \alpha}} a_{j}^{(\alpha)} \int_{\substack{[0,b] \\ d}} u^{2kz+h+j} du \quad (4.4)$$

$$= \sum_{\alpha} \sum_{\substack{j \mid \leq N \\ j \mid \leq N \\ \alpha}} a_{j}^{(\alpha)} \frac{d}{\prod} \frac{b^{2kpz+hp+jp+1}}{(2k_pz+hp+jp+1)}$$

Hence the function $3_1(z)$ can be analytically continued as claimed, with poles at the rational numbers (where h, k vary with α and $1 \le p \le d$)

$$Z = -\frac{h_p + j_p + 1}{2k_p} \qquad (4.5)$$

Then clearly the learning coefficient is

$$\lambda = \min \min_{\substack{\alpha \ j \le p \le d}} \left(\frac{h_p + l}{2k_p} \right). \qquad (J.l)$$

and its order (the power of $z + \lambda$ appearing in (4.4)), is

$$m = \max \# \{ p \mid \lambda = \frac{hp^{+}}{2kp} \}. \qquad (J.2)$$

From this we can derive (see [W,p.32, p.33] for the short version) e.g.

$$F_n^{\circ} \cong \lambda \log n - (m-1) \log \log n + F^R(\mathbf{F})$$
. (5.3).

RemarkAssuming $\mathcal{Y} > O$ on $\{W \mid K(w) = O\}$ we may in the aboveapply [W, Thm 2.3] so that h depends only on (W, K),and then multiply by \mathcal{Y} and still satisfy the conditions of[W, Thm 6.5(3)]. That is, we may awange on each \mathcal{A} for $h \in \mathbb{N}^d$ to be independent of the prior. Since k isclearly independent, we have that $\lambda = \lambda(p, q, W)$ is independent of \mathcal{Y} and agrees with the RLCT of $[W, Deg^2 2.7]$. Obrewe that the RLCT is defined as

$$RLCT = \inf \min_{w \in W} \left(\frac{h_p + 1}{k_p} \right).$$

Actually we only really need $\mathcal{G} > \mathcal{O}$ on those Mx for which the minimization in (5.1) is obtained. (check)

Remark The prior does affect the Schwartz distributions Dkm(4) on [W, p. 32]and hence the random variable \overline{S} in (5.3). <u>Def</u>^m We say $w \in W$ is <u>essential</u> for the tuple (P, Q, W) if the local resolution of singularities around P = w in $[W, Def^2.7]$ produces

$$\frac{h_p + 1}{k_p} = \lambda \quad (\text{the learning coefficient})$$

Necessarily $w \in W_0 = \{w \mid K(w) = 0\}$ is a singular point of W_0 , and it is an instance of the "worst" singularity type.

Local lambda

$$W = \left\{ w \in W^{(R)} \mid \mathcal{T}_{1}(w) \mathcal{P}_{0}, \mathcal{T}_{2}(w) \mathcal{P}_{0}, \dots, \mathcal{T}_{k}(w) \mathcal{P}_{0} \right\}$$

where $\pi_{1,...,\pi_{k}}: W^{(k)} \longrightarrow \mathbb{R}$ are real analytic and $W^{(k)} \in \mathbb{R}^{d}$ is open. This hypothesis is used in a subtle but crucial way in [W, Remark 2.14], i.e. the construction of the Ma, since we assume the resolution $g: V \rightarrow W$ makes $\pi_{i}(g(u)),...,\pi_{k}(g(u))$ normal crossing, i.e. (see [W, p.70]) for any $u_{0} \in U$, there is a local coordinate u_{i} s.t.

$$\pi_i(g(u)) = a_i(u) u^{k_{i1}} \cdots u^{k_{id}}$$

where $a_i(u)$ is everywhere nonzero. This means we may shrink the $O_p(b)$ above if necessary, so that on each M& the function $\pi_i(g(u))$ takes a consistent sign.



Now let $C \subseteq W$ be compact, defined by

$$C = \{ \omega \in W \mid f_1(\omega) \gg 0, \dots, f_l(\omega) \gg 0 \}$$

where f_j are real-analytic. Then adding the f_j to the list of functions to be "resolved" (i.e. made normal clossings) we may assume that the local wordinates u and $M_X \subseteq \mathcal{M}$ used above are adapted to Cin the sense that the boundary of $g^{-1}(C)$ is contained in $u_1 \cdots u_d = O$ in every such wordinate chart [W, Remark 2.12]. We assume <u>(i.e. realisability</u> $C_o = \{w \in K \mid K(w) = O\}$ is nonempty, and g' > O on C_o . for C

Lemma Let
$$\lambda$$
 denote the learning coefficient of the tuple $(P, 2, J, W)$
and λ_c the learning coefficient of $(P, 2, \mathcal{P}', C)$ for any
prior \mathcal{P}' as above. Then $\lambda_c \ge \lambda$.

Proof We have

$$g^{-1}(C) = \bigcup_{\alpha} g^{-1}(C) \cap M_{\alpha}.$$

In the local coordinates u appropriate to Ma, g⁻¹(C) Ma is given by f, (glu)) >0, ..., fr (g(u)) >0 and we may assume the balls Op(b) are chosen small enough so that for each index a in the relevant local coordinate

$$f_{j}(g(u)) = c_{j}(u)u^{g_{j1}} \dots u^{g_{jd}}$$

where we may assume cj(u) takes a fixed sign on Md. By construction Mx is a pwduct of (-6,0] and [0,6] intervals, and let mx be the number of occurrences of (-6,0]. We may assume these are the coordinates U1,..., Uma

Then
$$f_{i}(g(u)) \gtrsim 0, ..., f_{i}(g(u)) \gtrsim 0$$
 on Md
 $\iff c_{j}(u) u^{q_{j1}} \dots u^{q_{jd}} \gtrsim 0$ on Md $\forall i \leq j \leq d$
 $\iff (-i)^{\mathcal{X}_{j}} u^{q_{j}} \gtrsim 0$ on Md $\forall i \leq j \leq d$ $(\mathcal{T}_{j} = sign of (q_{j} = q_{ji} + \dots + q_{jd}))$
 $\iff (-i)^{\mathcal{Y}_{j}} u^{(q_{j})_{2}} \dots u_{md}^{(q_{j})_{md}} \geqslant 0$ on Md $\forall i \leq j \leq d$
 $\iff \mathcal{T}_{j} + (q_{j})_{1} + \dots + (q_{j})_{md}$ is even, $\forall i \leq j \leq d$.
So either $\mathcal{T}_{j} + (q_{j})_{1} + \dots + (q_{j})_{md}$ is even for all $1 \leq j \leq d$.
So either $\mathcal{T}_{j} + (q_{j})_{1} + \dots + (q_{j})_{md}$ is odd for some $1 \leq j \leq d$
in which case
 $g^{-1}(C) \cap Md = Md$, or this sum is odd for some $1 \leq j \leq d$
in which case
 $g^{-1}(C) \cap Md = \{u \in Md \mid u_{1} = \dots = u_{md} = 0\}$
Let Λ denote the set of α indices, $\Lambda' \leq \Lambda$ the set of indices for which
 $\mathcal{T}_{i} + (q_{i}) = 1 + (q_{i})_{i} + \dots + (q_{i})_{i}$

$$\partial_{j} + (Q_{j})_{1} + \cdots + (Q_{j})_{ma}$$
 is even.

$$\zeta_{c}(z) = \int_{c} K(w)^{z} \mathcal{I}'(w) dw$$

By restricting to W_{ε} , C_{ε} if necessary we may assume f > 0 on W and f' > 0 on C. Then in the partition of [W, Thm 6.5] we can assume $|g'(u)| = \phi^{pre}(u)u^{h}$ and $f(g(u))|g'(u)| = f(g(u))\phi^{pre}(u)u^{h}$



so that $\phi(u) = \mathcal{Y}(g(u)) \phi^{pre}(u)$, hence in (2.3)

$$\phi_{\alpha}^{\mathsf{t}}(u) = \mathcal{B}_{\alpha}(u)\phi(u) = \mathcal{B}_{\alpha}(u)\gamma(g(u))\phi^{\mathsf{pre}}(u).$$

We now compute for Re(z) > 0 using (2.3)

$$\begin{split} \zeta_{c}(z) &= \int_{c} K(w)^{z} \mathcal{G}'(w) dw \qquad (k,h \text{ depending on } d) \\ &= \sum_{\alpha} \int_{q^{-1}(c) Q M d} K(g(u))^{z} \mathcal{G}'(g(u)) \phi^{pre}(u) u^{h} du \end{split}$$

$$= \sum_{\alpha} \int_{\overline{g}'(c) \cap M_{\alpha}} u^{2kz} u^{h} \left[g'(g(u)) \phi^{pre}(u) \right] du$$

By the earlier argument (e.g. in (4.2)) we may ignore, for the purpose of computing A_c , the term in brackets involving f'. So we must compute the largest pole of

$$\sum_{\alpha} \int g^{-1}(c) \cap M_{\alpha}$$
 u du

$$= \sum_{\alpha \in \Lambda'} \int_{M_{\alpha}} u^{2kz+h} du + \sum_{\alpha \notin \Lambda'} \int_{\{u_{1}=\dots=u_{m_{\alpha}}=0\}}^{u^{2kz+h}} du$$

$$= \sum_{\alpha \in \Lambda'} \prod_{p=1}^{b} \int_{0}^{b} \frac{2k_{p}z+h_{p}}{du_{p}} du_{p}$$

$$= \sum_{\alpha \in \Lambda'} \frac{d}{p=1} \frac{b^{2k_{p}z+h_{p}}}{2k_{p}z+h_{p}+1} \qquad (9.1)$$

Note that these tuples k, $h \in \mathbb{N}^d$ are the same as used to define λ , so that

$$\lambda_{c} = \min \min \left(\frac{h_{p} + |}{zk_{p}} \right)$$

$$\neq \min \min \left(\frac{h_{p} + |}{zk_{p}} \right) = \lambda.$$

$$(10.1)$$

$$\neq \min \left(\frac{h_{p} + |}{zk_{p}} \right) = \lambda.$$

Resolutions and the posterior



If we set $x = x_1y_1$ and $y = y_1$ this becomes a standard Gaussian on x_1y_70

$$e^{-(1+\frac{x^2}{y^2})y^2} = e^{-x^2-y^2}$$
(10.3)

Put differently, if $g: [0,b]^2 \longrightarrow \mathbb{R}^2$ is $g(x_1,y_1) = (x_1y_1,y_1)$ and $f: \mathbb{R}^2 \longrightarrow \mathbb{R}$ is $f(x_1y) = \exp(-x^2 - y^2)$ then $f(g(x_1,y_1)) = \exp(-(1+x_1^2)y_1^2)$. This is (one wordinate patch of) the resolution of $x^2 + y^2$ [W, Example 2.7]. The resolution procedure [W, I.4.2] expresses the measure exp $(-n\beta K_n(w)) S(w) dw$, where normalisation is the Bayesian posterior, as asymptotically approximated by a sum

$$\sum_{\alpha} \exp\left(-n\beta u^{2k} + \int n\beta u^{k} \mathcal{F}_{n}(u)\right) \phi^{*}(u) |u^{h}| du \quad (11.1)$$

For the sake of building intuition let us ignore the stochastic process $\overline{S}_n(u)$, and the $\phi^*(u)$ (which is made up of the partition of unity and prior). The $|u^h|$ is a Jacobian factor, which corresponds to pushing forward the measure

$$\exp\left(-n\beta u^{2k}\right) du = \exp\left(-n\beta u_{1}^{2k} - - u_{d}^{2kd}\right) \quad (11.2)$$

along $g: U \longrightarrow W$. Recall that the index α is associated with a product $[0, b]^d$ in some local coordinates $[W, \operatorname{Remark} 2.14]$ and (11.1) is to be interpreted as each of these local patches $M\alpha \equiv [0, b]^d$ contributing additively to approximate the posterior.

So what kind of function is (11.2) on $M_{x} = [0, b]^{d}$?





Similarly for k = (1,0). The other cases are similar, but the approach to the coordinate hyperplanes is more gradual than Gaussian, e.g. in k = (0,2) we have $exp(-n\beta y_2^4)$ so the measure is more spread out along the u_2 axis.

In higher dimensions, e.g. d = 3, the picture is similar. We pick some subject of the coordinate hyperplanes (those with the corresponding coordinate of knonzew), place a large measure along those hyperplanes and "round off the corners" near u = 0. Higher indices k_i correspond to flatter approaches along that coordinate.



We now imagine a collection of such measures on sets Md, each with their own map $g_{\alpha}: M_{\alpha} \longrightarrow V$ to some neighborhood V of W, such that a weighted sum over a of the pushforward of these measures along each g_{α} approximates the posterior on V. To borow from [W, Fig 2.5, Fig 2.10]



This suggests a potential generalisation of implicit variational inference, which we term <u>microscopic variational inference</u> (MVI). We fix an integer S and let Go.,..., Gos be feedforward networks with weight vectors Oi; whose outputs are in IR^d and whose inputs are in $[0,1]^d$. For each $1 \le c \le S$ we choose a vector $k_i = (k_{i1}, ..., k_{id}) \in IN^d$ and we introduce additional weights $a_{i1}..., a_{id}$ which are logits for weighting each network.

We want to sample, for each i, from a normalisation of $exp(-t_i u^{k_i})$ where t_i is some (inverse) temperature, pass that sample through Go: and sum these outputs with weights from a softmax of $(a_1, ..., a_s)$. Set $Z_i := \int_{[0,1]}^{\infty} d \exp(-t_i u^{k_i}) du$. We assume some <u>target weight</u> vector $P \in W$ is specified (as in (13.1)) and we wish to approximate the posterior in a small neighborhood of P, e.g. to estimate the local RLCT. So we use the output of each Go_i as a <u>delta</u> from P, so that our probability distribution associated to the parameters $(O_1, \dots, O_s, O_1, \dots, O_s)$ is

$$\sum_{i=1}^{S} \frac{\exp(a_i)}{\sum_{j} \exp(a_j)} (G_{0_i})_* (\frac{1}{Z_i} \exp(-t_i u^{k_i})) + P \quad (14.1)$$

where (Goi) means pushforward along the function computed by the network Goi.

The idea behind MVI is that in the singular case, it is not enough to think about modes and Gaussians; you need to think in terms of the distributions $exp(-n\beta u^R)$ on $p(\overline{U})$, \overline{U} more generally, and how to sum transformations of these. By universality and the resolution theorem (eliding the difference between partition of unity and the softmax) for <u>some</u> weight $O_{1,...,}O_{s}$ you can approximate the three posterior this way (asymptotically, as in [W, p. 33], but this enough to compute local RLCTs) as $S \rightarrow \infty$, the ki's range over all tuples and the depths of the Go; go to infinity.

Remark Consider the following graphs



 (\overline{S})

For x = 0 we have

$$\frac{d}{dx} \exp(-x^{2}) = -2x \exp(-x^{2}) \approx -2x + O(x^{3})$$

$$\frac{d}{dx} \exp(-x^{3}) = -3x^{2} \exp(-x^{2}) \approx -3x^{2} + O(x^{4})$$
(15.2)

So the larger k_i is, the flatter $exp(-n\beta u^k)$ is as we approach the ith coordinate hyperplane in $[0, b]^d \cong Md$. But it is actually the transformed measure (11.1) we are interested in It is not simple to explain the interaction of k, h in this formula directly, so to understand how $h_i k$ combine to cletermine the local behaviour of the posterior, we switch to talking about volumes.

Consider a regular model in a neighborhood of the twe parameter, which we may assume is $\omega_0 = 0$, so that near ω_0 (possibly atter changing wordinates)

$$K(\omega) \approx \sum_{i=1}^{d} \omega_i^2 \qquad (16.1)$$

Consider the volume of almost correct parameters

$$V(t) = \int_{K(\omega) < t}^{\prime} \mathcal{I}(\omega) d\omega . \qquad (16.2)$$

We know the prior is (almost) irrelevant for the purposes of understanding the RLCT, so for simplicity we take it to be uniform so that V(t) is the volume of all d-ball of radius $t^{1/2}$

$$V(t) \propto t^{a/2} \tag{16.3}$$

Hence for
$$a > 0$$
, $a \neq 1$ we have

$$\log \left\{ \frac{\sqrt{at}}{\sqrt{t}} \right\} = \log \left\{ \frac{a^{d/2} t^{d/2}}{t^{d/2}} \right\} \quad (16.4)$$

$$= \log \left\{ a^{d/2} \right\} = \frac{d}{2} \log a.$$
Hence [W, Theorem 7.1] gives $\lambda = \frac{d}{2}$, as the RLCT.

In the regular cone 2λ is the number of parameter. We now proceed to justify the claim that 2λ is an appropriate <u>definition</u> of the effective number of parameters ingeneral, starting with a mildly singular case.

sd14) (16)

$$\frac{dT}{dt}$$
Remark In the case $d = 2$ the relevant diagram is
$$\frac{dT}{dt} = \frac{1}{2} \operatorname{The relevant} \operatorname{diagram is}$$

$$\frac{dT}{dt} = \frac{1}{2} \operatorname{The relevant} \operatorname$$

Note that, at least near w_0 , $\{w | K|w\}=0\}$ has a free vaniable w_i for i > d'and hence dimension d-d', so the normal bundle is d' dimensional. Hence $2\lambda = d'$ is the number of directions we can vary w_0 which change the KL divergence. That is, 2λ is the effective number of parameters in the model near w_0 .

In the case d = 3, d' = 2 the relevant cliagram is (supposing $K(\omega) \approx \omega_1^2 + \omega_2^2$)



<u>General singular case</u> in general while $W_o = \{w \mid K(w) = 0\}$ is a real analytic set it is not a submanifold, and near a singular point there is no naive notion of a "normal direction". So we cannot simply count the number of normal directions (for the same reason the tubular neighborhood theorem fails in algebraic geometry) as we have done in (17.1) (where every direction was normal, as W_o was a point) and (18.1) (where W_o was a line and there are two normal directions).

However the volume ratios, which give the dimension of the normal bundle in the regular and "mildly singular" cases; are still well-defined and by [W, Tum 7.1]

$$\lambda = \lim_{t \to 0} \frac{\log\{\sqrt{(at)}/\sqrt{(t)}\}}{\log a} \qquad (a>0, a\neq 1)$$

Here is how to think about this: suppose someone gave you the loss surface $\{(\omega, \kappa(\omega)) | \omega \in W\}$ and a point $\omega_0 \in W$. You start computing volumes V(0.99t), V(t) for some small values of t. This involves cutoffs and priors but the ratio V(0.99t)/V(t) is basically independent of these. And that someone says to you: if this were a mildly singular model, what would the dimension of the normal bundle be?

That is, what dimensionality would these volumes have, if they were d'balls for some d'? The answer to this question is what 22 computes, hence why Watanabe refers to this as volume *elimension* in [W, Remark 7.1].

Of course 22 need not be an integer, but this corresponds to the fact that there is no "true" normal bundle, so there cannot really be an integral "effective" number of parameters in a singular model at a degenerate critical point.

<u>Remark</u> The RLCT is not a measure of "flatness" in the same sense as e.g. the determinant of the Hessian : note that in the regular case it is d/2regardless of the eigenvalues. Instead the RLCT is a kind of "count" of the number of flat clirections, at least in directly, since

small $\lambda \implies$ small normal bunche \implies many flat directions.

5d14

Learning coefficient vs. generalisation error

Suppose as above we have (P, Q, W, Y) and $C \subseteq W$ compact, with q(x) realisable by some parameter in C, let λ_c be the learning coefficient for C. Then with (P, Q, C, Y') as our input to $[W, \S6]$ we have $(\alpha t \beta = 1)$

$$\mathbb{E}\left[B_{g}^{c}\right] = \frac{\lambda_{c}}{n} + o\left(\frac{1}{n}\right) \qquad (10.1)$$

$$B_g^{C} := \mathbb{E}_{\times} \left[\log \frac{q(x)}{\mathbb{E}_{\omega}} \left[p(X|\omega) \right] \right]$$

where \mathbb{F}_{X} means $\int q(x)(-) dx$ and

$$\mathbb{E}_{w}\left[p(X|w)\right] = \int p(w|D_{n}) p(X|w) dw$$
$$= \frac{1}{Z_{n}^{o}} \int p(X|w) \exp\left(-n\beta K_{n}(w)\right) dw$$

is the <u>C-restricted predictive distribution</u>. So Bg is the KL divergence between the true distribution and the C-restricted predictive distribution, and the stochasticity involved in sampling Dn makes Bg a random variable. The expectation in (20.1) integrates out this stochasticity.

Conclusion : the big picture

So the picture is the following: supposing q is realisable by (p, W), the set $W_0 = \{w \in W \mid K(u) = 0\}$ is nonempty, and in general a complex real analytic set with multiple components and singularities:



(21)

Shown are three points $w_0^1, w_0^2, w_0^3, w_0^4 \in W_0$ and compact $C_i \subseteq W$ containing them. The point w_0^4 is an isolated zero of K(w) and while it may nonetheless be a degenerate critical point, let us imagine it to be nondegenerate so $\lambda_{c_4} = d/2$. The point w_0^3 is a smooth point of Wo and looks perhaps like the "mild" case above, i.e. locally $K(w) = \sum_{i=1}^{d} w_0^3$, with $\lambda_{c_3} = d'/_2 < d/2$, and locally d'honest normal directions to the set of the parameters.

The point ω_0^2 is a more complex singularity, and ω_0' more complex still, so that we have

$\lambda_{c_1} < \lambda_{c_2} < \lambda_{c_3} < \lambda_{c_4} = d/2$

Hence at a given dataset size, we have the same relationship between The Bayesian generalisation errors. The models near wo' have the least number of effective parameters, or degrees of freedom, and generalise best.