# Statistical Learning Theory of Quasi-Regular Cases

Koshi Yamada*      Sumio Watanabe*

October 1, 2018

## Abstract

Many learning machines such as normal mixtures and layered neural networks are not regular but singular statistical models, because the map from a parameter to a probability distribution is not one-to-one. The conventional statistical asymptotic theory can not be applied to such learning machines because the likelihood function can not be approximated by any normal distribution. Recently, new statistical theory has been established based on algebraic geometry and it was clarified that the generalization and training errors are determined by two birational invariants, the real log canonical threshold and the singular fluctuation. However, their concrete values are left unknown. In the present paper, we propose a new concept, a quasi-regular case in statistical learning theory. A quasi-regular case is not a regular case but a singular case, however, it has the same property as a regular case. In fact, we prove that, in a quasi-regular case, two birational invariants are equal to each other, resulting that the symmetry of the generalization and training errors holds. Moreover, the concrete values of two birational invariants are explicitly obtained, the quasi-regular case is useful to study statistical learning theory.

## 1   Introduction

A lot of statistical learning machines which are being applied to pattern recognition, bioinformatics, robotic control, and artificial intelligence have hidden variables, hierarchical layers, and submodules, because they are used to estimate the structure of the true distributions. In such learning machines, the map taking parameters to probability distributions is not one-to-one and the Fisher information matrices are singular, hence they are called singular learning machines. For example, three-layered neural networks, normal mixtures, hidden Markov models, Bayesian networks, and reduced rank regressions are singular learning machines [1, 2, 4, 5, 6, 10]. If a statistical model is singular, then either

*Department of Computational Intelligence and Systems Science, Tokyo Institute of Technology, Mail box:G5-19, 4259 Nagatsuta, Midori-ku, Yokohama, 226-8502, Japan E-mail:yamada.k.am@m.titech.ac.jp,swatanab@dis.titech.ac.jp

the maximum likelihood estimator is not subject to the normal distribution even asymptotically or the Bayes posterior distribution can not be approximated by any normal distribution. Hence it has been difficult to study their learning performace and to estimate the generalization error from the training error.

Recently, new statistical theory has been established based on algebraic geometrical method [11, 15, 13, 14] and it was clarified that the generalization and training errors in Bayes estimation, $G_n$ and $T_n$, are given by two birational invariants, the real log canonical threshold $\lambda$ and singular fluctuation $\nu$ by the formulas,

$$\mathbb{E}[G_n] = \Big(\frac{\lambda - \nu}{\beta} + \nu\Big)\frac{1}{n} + o(\frac{1}{n}), \tag{1}$$

$$\mathbb{E}[T_n] = \Big(\frac{\lambda - \nu}{\beta} - \nu\Big)\frac{1}{n} + o(\frac{1}{n}), \tag{2}$$

where $\mathbb{E}[\ ]$ shows the expectation value over all training sets, $n$ is the number of training samples and $\beta$ is the inverse temperature of the Bayes posterior distribution. Based on this relation, we can define an information criterion which enables us to estimate the generalization error from the training error [13].

It is well known that, if the true distribution and the statistical model are in a regular case, then $\lambda = \nu = d/2$ holds where $d$ is the dimension of the parameter space. In this case, the symmetry of the generalization and training errors holds,

$$\mathbb{E}[G_n] = \frac{d}{2n} + o(\frac{1}{n}), \tag{3}$$

$$\mathbb{E}[T_n] = -\frac{d}{2n} + o(\frac{1}{n}), \tag{4}$$

for arbitrary $0 < \beta \leq \infty$. This case corresponds to the well-known Akaike Information criterion for regular statistical models. However, if they are not in a regular case, neither of them is equal to $d/2$ in general. Therefore, in order to study singular learning machines, researches on two birational invariants are necessary.

In the present paper, in order to investigate the mathematical structure of birational invariants, we firstly introduce a new concept, a quasi-regular case, which satisfies the relation,

Regular $\subsetneq$ Quasi-Regular $\subsetneq$ Singular.

In other words, a quasi-regular case is not a regular case, however, it has the same properties as the regular case. In fact, we prove that, in quasi-regular cases, both birational invariants are equal to each other, $\lambda = \nu$, and the symmetry of the generalization and training errors holds. In a quasi-regular case, two birational invariants are obtained explicitly, hence it is a useful concept in researches of statistical learning theory.

2

# 2 Framework of Bayes Learning

In this section, we summarize the framework of the Bayes learning, and introduce the well-known results.

## 2.1 Generalization and Training Errors

Firstly, we define the generalization and training errors. Let $N$, $n$ and $d$ be natural numbers. Let $X_1, X_2, ..., X_n$ be random variables on $\mathbb{R}^N$ which are independently subject to the same probability density function as $q(x)$. Let $p(x|w)$ be a probability density function of $x$ for a parameter $w \in W \subset \mathbb{R}^d$, where $W$ is a set of parameters. The prior distribution is represented by the probability density function $\varphi(w)$ on $W$. For a given training set

$$X^n = \{X_1, X_2, ..., X_n\},$$

the posterior distribution is defined by

$$p(w|X^n) = \frac{1}{Z_n} \prod_{i=1}^{n} p(X_i|w)^\beta \varphi(w) dw,$$

where $0 < \beta < \infty$ is the inverse temperature and $Z_n$ is the normalizing constant. The case $\beta = 1$ is most important because it corresponds to the strict Bayes estimation. The expectation value over the posterior distribution is denoted by

$$\mathbb{E}_w[\ \ ] = \int (\ \ ) p(w|X^n) dw.$$

The predictive distribution is defined by

$$p(x|X^n) = \mathbb{E}_w[p(x|w)].$$

The generalization and training error, $G_n$ and $T_n$, are respectively defined by

$$
\begin{aligned}
G_n &= \int q(x) \log \frac{q(x)}{p(x|X^n)} dx, \\
T_n &= \frac{1}{n} \sum_{i=1}^{n} \log \frac{q(X_i)}{p(X_i|X^n)}.
\end{aligned}
$$

The generalization error shows the Kullback-Leibler distance from the true distribution to the estimated distribution. The smaller the generalization error is, the better the learning result is. However, we can not know the generalization error directly, because calculation of $G_n$ needs the expectation value over the unknown true distribution $q(x)$. On the other hand, the training error can be calculated using only training samples, in practice, as the log likelihood function. Hence one of the main purposes of statistical learning theory is to clarify the mathematical relation between them.

## 2.2 Two Birational Invariants

Secondly, we define two birational invariants.

The Kullback-Leibler distance from the true distribution $q(x)$ to a parametric model $p(x|w)$ is defined by

$$K(w) = \int q(x) \log \frac{q(x)}{p(x|w)} dx.$$

Then $K(w) = 0$ if and only if $q(x) = p(x|w)$. In this paper, we assume that there exists a parameter $w_0$ which satisfies $q(x) = p(x|w_0)$ and that $K(w)$ is an analytic function of $w$.

**Definition 1 (Real Log Canonical Threshold)** *The zeta function of statistical learning is defined by*

$$\zeta(z) = \int K(w)^z \varphi(w) dw.$$

*Then $\zeta(z)$ is a holomorphic function on the region $Re(z) > 0$, which can be analytically continued to the unique meromorphic function on the entire complex plane [11]. All poles of the zeta function are real, negative, and rational numbers. If its largest pole is $(-\lambda)$, then the real log canonical threshold is defined by $\lambda$. The order of the pole $z = -\lambda$ is referred to as a multiplicity $m$.*

**Definition 2 (Singular Fluctuation)** *The functional variance is defined by*

$$V_n = \sum_{i=1}^{n} \{ \mathbb{E}_w[(\log p(X_i|w))^2] - \mathbb{E}_w[\log p(X_i|w)]^2 \}.$$

*Then it was proved [13] that the expectation value*

$$\nu = \frac{\beta}{2} \lim_{n \to \infty} \mathbb{E}[V_n]$$

*exists. The constant $\nu$ is called the singular fluctuation.*

**Theorem 1** *The expectation values of the generalization and training errors are given by eq.(1) and eq.(2). Therefore*

$$\mathbb{E}[G_n] = \mathbb{E}[T_n] + \frac{2\nu}{n} + o(\frac{1}{n}).$$

(Proof) This theorem was proved in [15, 13]. (Q.E.D.)

**Remarks**. (1) The real log canonical threshold and the singular fluctuation are invariant under a birational transform

$$\begin{aligned} w &= g(w'), \\ p(x|w) &\mapsto p(x|g(w')), \\ \varphi(w) &\mapsto \varphi(g(w'))|g'(w')|, \end{aligned}$$

4

where $|g'(w')|$ is the Jacobian determinant. Such constants are called birational invariants.

(2) The real log canonical thresholds for several learning machines were clarified [1, 2, 16] using resolution of singularities. However, the singular fluctuation has been left unknown. This paper provides the first result which clarifies the concrete values of singular fluctuation in a singular case.

(3) The real log canonical threshold is a well known birational invariant in algberaic geometry, which plays an important role in higher dimensional algberaic geometry. The singular fluctuation was found in statistical learning theory.

## 2.3 Regular and Singular

Thirdly, we define regular and singular cases.

**Definition 3** *A pair of the true distribution and the parametric model, $(q(x), p(x|w))$ is called to be in a regular case if and only if the set $\{w; q(x) = p(x|w)\}$ consists of a single element $w_0$ and Fisher information matrix*

$$\int \nabla \log p(x|w_0)(\nabla \log p(x|w_0))^T q(x)dx$$

*is positive definite. Otherwise, it is called to be in a singular case.*

For a regular case, the real log canonical threshold and the singular fluctuation have been completely clarified.

**Theorem 2** *If a pair $(q(x), p(x|w))$ is in a regular case, then $\lambda = \nu = d/2$, where d is the dimension of the parameter.*

(Proof) This theorem was proved in [15]. (Q.E.D.)

# 3 Main Results

In this section, we define a quasi-regular case. This concept is firstly proposed by the present paper. Also the main theorem is introduced.

**Definition 4** *¿**Quasi-Regular Case**j. Assume that there exists a parameter $w_0 \in W^o$ such that $q(x) = p(x|w_0)$. Without loss of generality, we can assume that $w_0$ is the origin $w_0 = 0$. The original parameter is denoted by $w = (w_1, w_2, ..., w_d)$. Let g and $\Delta d_1, \Delta d_2, ..., \Delta d_g$ be natural numbers which satisfy*

$$\Delta d_1 + \Delta d_2 + \cdots + \Delta d_g = d$$

*and $\Delta d_0 = 0$. We define*

$$d_j = \Delta d_0 + \cdots + \Delta d_j \ \ (j = 0, \cdots, g)$$

*and a function $u = (u_1, u_2, ..., u_g) \in \mathbb{R}^g$ of the paramater $w \in \mathbb{R}^d$ by*

$$
\begin{aligned}
u_1 &= \prod_{j=1}^{d_1} w_j, \\
u_2 &= \prod_{j=d_1+1}^{d_2} w_j, \\
\cdots &= \cdots, \\
u_g &= \prod_{j=d_{g-1}+1}^{d_g} w_j.
\end{aligned}
$$

*If there exist constants $c_1, c_2 > 0$ such that, for arbitrary $w \in W$,*

$$
c_1(u_1^2 + \cdots + u_g^2) \leq K(w) \leq c_2(u_1^2 + \cdots + u_g^2),
$$

*then the pair $(q(x), p(x|w))$ is called to be in a quasi-regular case.*

**Remark.** (1) If $g = d$, then

$$
\{w; q(x) = p(x|w)\} = \{0\}
$$

and the quasi-regular case corresponds to the regular case. Hence a quasi-regular case contains a regular case as a special one.

(2) If $d \neq g$, then "$K(w) = 0 \Longleftrightarrow w = 0$" does not hold, because, for at least one variable $w_j$, $K(0, 0, .., w_j, 0, .., 0) = 0$. Hence a quasi-regular case with $d \neq g$ is not a regular case but a singular case.

(3) There are singular cases which are not contained in quasi-regular cases. Therefore,

$$
\text{Regular} \subsetneq \text{Quasi-Regular} \subsetneq \text{Singular}
$$

holds. The present paper shows in Theorem 3 that a quasi-regular case is not a regular case, however, it has the same property as a regular case.

**Example.1** Let a statistical model be

$$
p(x, y|w) = \frac{r(x)}{\sqrt{2\pi}} \exp(-\frac{1}{2}(y - ax^2 - b\tanh(cx))^2),
$$

where $w = (a, b, c)$ is the parameter and $r(x)$ is the probability density function of $x$. If the true distribution is given by $q(x, y) = p(x, y|0, 0, 0)$, then by using

$$
u_1 = a, \quad u_2 = bc,
$$

it follows that

$$
K(w) = \frac{1}{2} \int (ax^2 + b\tanh(cx))^2 r(x) dx
$$

satisfies the condition for a quasi-regular case with $g = 2$, because $x^2$ and $\tanh(cx)/c$ is linearly independent. In fact there exist $c_1, c_2 > 0$ such that

$$c_1(a^2 + (bc)^2) \leq K(w) \leq c_2(a^2 + (bc)^2).$$

Hence the set of true parameters consists of the union of two lines,

$$\{w; q(x,y) = p(x,y|w)\} = \{a = 0, bc = 0\}.$$

**Example.2** Let a statistical model be

$$p(x,y|w) = \frac{r(x)}{\sqrt{2\pi}} \exp(-\frac{1}{2}(y - ax - b\tanh(cx))^2),$$

where $w = (a,b,c)$ is the parameter and the true distribution is given by $q(x,y) = p(x,y|0,0,0)$. Then because $x$ and $\tanh(cx)/c$ is not linearly independent as $c \to 0$, hence this case does not satisfies the quasi-regular condition. In this case

$$c_1((a + bc)^2 + b^2c^6) \leq K(w) \leq c_2((a + bc)^2 + b^2c^6).$$

Example.2 resembles Example.1, however, from the viewpoint of statistical learning theory, they are different.

**Example.3** Let a statistical model be

$$p(x,y,z|w) = \frac{r(x,y)}{\sqrt{2\pi}} \exp(-\frac{1}{2}(z - f(x,y,w))^2),$$

where

$$
\begin{aligned}
f(x,y,w) &= a_1\sin(b_1x) + a_2x\sin(b_2x) \\
&\quad + a_3\sin(b_3y) + a_4y\sin(b_4y),
\end{aligned}
$$

and $w = \{(a_i, b_i)\}$ is the parameter and the true distribution is given by $q(x,y,z) = p(x,y,z|0)$. Then $(q(x,y,z), p(x,y,z|w))$ is in a quasi-regular case with $g = 4$.

The following is the main theorem of the present paper.

**Theorem 3 (Main Theorem).** *Assume that the pair $(q(x), p(x|w))$ is in a quasi-regular case and that $\varphi(w) > 0$ on $W$. Then the real log canonical threshold and the singular fluctuation are given by*

$$\lambda = \nu = \frac{g}{2}$$

*and*

$$m = d - g + 1.$$

**Corollary 1** *Assume that the pair $(q(x), p(x|w))$ is in a quasi-regular case and that $\varphi(w) > 0$ on $W$. For arbitrary $0 < \beta < \infty$ the symmetry of the generalization and training errors holds,*

$$\mathbb{E}[G_n] = \frac{g}{2n} + o(\frac{1}{n}),$$
$$\mathbb{E}[T_n] = -\frac{g}{2n} + o(\frac{1}{n}).$$

**Remarks.**(1) The above theorem shows the generalization and training errors for Bayes estimation. In the quasi-regular case, they have the same property as those in regular cases, however, the generalization and training errors of the maximum likelihood estimation is different from regular case in general.
(2) In the maximum likelihood method, the training error of a singular case is far smaller than that of a regular case, whereas the generalization error of a singular case is far larger than that of a regular case. From the viewpoint of the maximum likelihood method, the quasi-regualr case is contained in the singular case. In the present paper, we prove that the quasi-regular case has the same property as the regular case from the viewpoint of the Bayes estimation.

## 4    Proofs

In this section, we prove the main theorem. At first, we derive the real log canonical threshold of the quasi-regular case.

**Lemma 1** *The real log canonical threshold and its order are given by $\lambda = g/2$ and $m = d - g + 1$ respectively.*

(Proof) Since each function $\{u_j; j = 1, 2, ..., g\}$ does not have common variable $w_k$, the real log canonical threshold is given by the sum of individual real log canonical thresholds (Remark 7.2 in [15]) defined by

$$\zeta_j(z) = \int \prod_{i=d_{j-1}+1}^{d_j} (w_i)^{2z} dw_i$$
$$= \frac{C}{(z + 1/2)^{d_j - d_{j-1}}} + \cdots + .$$

Hence $\lambda$ is equal to $g$ times $1/2$, hence $\lambda = g/2$. The multiplicity is also given by

$$m = d_1 + d_2 - d_1 + \cdots + d_g - d_{g-1} - (g - 1)$$
$$= d - g + 1,$$

which shows the Lemma. (Q.E.D.)

**Definition 5** *For a given pair of the true distribution $q(x)$ and the parametric model $p(x|w)$, the log density ratio function is defined by*

$$f(x, w) = \log \frac{q(x)}{p(x|w)}.$$

The following lemma shows that the log density ratio function of the quasi-regular case is represented by $g$ linearly independent functions.

**Lemma 2** *Assume that the pair $(q(x), p(x|w))$ is in a quasi-regular case. Then there exists a set of functions $\{e_j(x, u); j = 1, 2, ..., g\}$ which are analytic functions of $u$ and*

$$f(x, w) = \sum_{j=1}^{g} u_j e_j(x, u)$$

*in an open neighborhood of $u = 0$.*

(Proof) Let us define a function

$$F(t) = t + e^{-t} - 1.$$

for $t \in \mathbb{R}^1$. Then $F(0) = 0$, $F'(0) = 0$, and $F''(0) = 1$, resulting that $F(t) \geq 0$ and that $F(t) = 0$ if and only if $t = 0$. Moreover, $F(t) \cong (1/2)t^2$ for small $|t|$. Therefore,

$$
\begin{aligned}
K(w) &= \int q(x) F\left(\log \frac{q(x)}{p(x|w)}\right) dx \\
&= \int q(x) F(f(x, w)) dx \\
&\cong \frac{1}{2} \int q(x) f(x, w)^2 dx.
\end{aligned}
\tag{5}
$$

By the assumption of the quasi-regular case, $K(w) = 0$ if and only if $u_1 = u_2 = \cdots = u_g = 0$, which is equivalent to $f(x, w) \equiv 0$. That is to say, $f(x, w)$ is contained in the ideal of analytic functions generated by $u_1, u_2, ..., u_g$. Hence there exist a set $\{e_j(x, u)\}$ of analytic functions of $u$, which satisfies

$$f(x, w) = \sum_{j=1}^{g} u_j e_j(x, u).$$

Therefore, we obtained the Lemma. (Q.E.D.)

In the following lemma, we show that the quasi-regular case has the generalized Fisher information matrix.

**Lemma 3** *The $g \times g$ matrix $I(u)$ is defined by*

$$I_{ij}(u) \equiv \int q(x) e_i(x, u) e_j(x, u) dx.$$

*Then $I(u)$ is positive definite in an open neighborhood of $u = 0$.*

9

(Proof) By Lemma 2 and eq.(5), in the neighborhood of $u = 0$,

$$K(w) = \frac{1}{2}(u \cdot I(u)u).$$

By the condition of the quasi-regular case,

$$c_1 \sum_{j=1}^{g} u_j^2 \leq K(w).$$

Hence the minimum eigenvalue of $I(u)$ is positive, which shows $I(u)$ is positive definite. (Q.E.D.)

The following definition and lemma show that the empirical loss function of the quasi-regular case has the same decomposition as that of the regular case.

**Definition 6** *A random process $\xi_n(u) \in \mathbb{R}^g$ is defined by*

$$\xi_n(u) = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} \{\frac{1}{2}I(u)u - e(X_i, u)\}$$

*where*

$$e(x, u) = (e_1(x, u), e_2(x, u), ..., e_g(x, u))^T.$$

**Lemma 4** *The empirical loss function defined by*

$$K_n(w) = \frac{1}{n} \sum_{i=1}^{n} f(X_i, w)$$

*is represented by*

$$K_n(w) = \frac{1}{2}(u, I(u)u) - \frac{1}{\sqrt{n}} u \cdot \xi_n(u)$$

*in the neighborhood of $u = 0$. Moreover, the random process $\xi_n(u)$ converges to the gaussian process $\xi(u)$ that satisfies*

$$\mathbb{E}[\xi(0) \cdot I(0)^{-1}\xi(0)] = g.$$

(Proof) The empirical loss function is given by

$$K_n(w) = K(w) - \frac{1}{n} \sum_{i=1}^{n} \{K(w) - f(X_i, w)\}.$$

By combining this equation with the definition of $\xi_n(u)$, the first half of the Lemma is obtained. For the second half, the convergence $\xi_n(u)$ is derived from the general empirical process theory. Moreover,

$$\mathbb{E}[\xi_n(0) \cdot I(0)^{-1}\xi_n(0)]$$
$$= \mathbb{E}[\mathrm{tr}(I(0)^{-1}\xi_n(0)\xi_n(0)^T)] = g,$$

where we used the covariance matrix of $\xi_n(0)$

$$\mathbb{E}[\xi_n(0)\xi_n(0)^T] = \int q(x)e(x,0)e(x,0)^T dx = I(0),$$

which completes the Lemma. (Q.E.D.)

In the quasi-regular case, the relation between $w = (w_1, w_2, ..., w_d)$ and $u = (u_1, u_2, ..., u_g)$ is important. The following lemma shows the property of the quasi-regular case. This lemma does not hold in general singular cases.

**Lemma 5** *When $n$ tends to infinity,*

$$\prod_{j=1}^{g} \delta\left(\frac{u_j}{\sqrt{n}} - \prod_{k=d_{j-1}+1}^{d_j} w_k\right) \cong c_3 (\log n)^{m-1} \prod_{j=1}^{d} \delta(w_j)$$

*where $m = d - g + 1$ and $c_3 > 0$ is a constant.*

(Proof) Firstly, we prove that the delta function with variables $\mathbf{x} = (x_1, x_2, ..., x_d)$ in $M \equiv [0,1]^d$
$$D(t, \mathbf{x}) = \delta(t - x_1 x_2 \cdots x_d)$$
has asymptotic expansion for $t \to 0$,

$$\begin{aligned}
D(t, \mathbf{x}) &= \frac{(-\log t)^{d-1}}{(d-1)!} \prod_{k=1}^{d} \delta(x_k) \\
&\quad + o((-\log t)^{d-2}).
\end{aligned} \tag{6}$$

Let $\phi(\mathbf{x})$ be an arbitrary $C^\infty$-class function of $\mathbf{x}$ whose support is containd in

$$D_t(\phi) \equiv \int_M D(t, \mathbf{x})\phi(\mathbf{x})d\mathbf{x}.$$

Then its Mellin transform is

$$\int D_t(\phi)t^z dt = \int_M \prod_{i=1}^{d}(x_i)^z \phi(\mathbf{x})d\mathbf{x},$$

where $M$ is the compact set that is the support of $\phi$. Without loss of generality By using Taylor expansion

$$\phi(\mathbf{x}) = \phi(0) + \mathbf{x} \cdot \nabla\phi(0) + \cdots,$$

we have the asymptotic expansion,

$$\int D_t(\phi)t^z dt = \frac{1}{(z+1)^d}\phi(0) + \cdots.$$

11

Therefore

$$\int D_t(t,\mathbf{x})t^z dt = \frac{1}{(z+1)^d}\prod_{k=1}^{d}\delta(x_k)+\cdots$$

for $\mathbf{x}\in[0,1]^d$. By using inverse Mellin transform, we obtained eq.(6). Secondly, let us prove the Lemma. By using eq.(6), for each $u_j$,

$$\delta(\frac{u_j}{\sqrt{n}}-\prod_{k=d_{j-1}+1}^{d_j}w_k)$$

$$\propto (\log n)^{d_j-d_{j-1}-1}\prod_{j=d_{j-1}+1}^{d_j}\delta(w_j)$$

when $n\to\infty$. By summing up these relations for $j=1,2,...,g$, Lemma is obtained. (Q.E.D.)

Let us return to the proof of the Main theorem.

**(Proof of Main Theorem)** It was proved by eq.(6.4) in [15] that the expectation value of $K_n(w)$ is given by two birational invariants,

$$\mathbb{E}[\mathbb{E}_w[K_n(w)]] = \frac{\lambda}{n\beta}-\frac{\nu}{n}+o(\frac{1}{n}).$$

Since we have already obtained the value of $\lambda$ in Lemma1, that is to say, $\lambda = g/2$, we can derive the value of $\nu$ by calculating $\mathbb{E}[\mathbb{E}_w[K_n(w)]]$. The posterior distribution is represented by the empirical loss function by

$$p(w|X^n)\propto\exp(-n\beta K_n(w))\varphi(w)dw.$$

The integration of the outside of the neighborhood of $u=0$ with respect to the posterior distribution goes to zero with the smaller order than $\exp(-\sqrt{n})$ as Lemma 6.3 in [15], hence we can restrict the integrated region to the neighborhood of $u=0$. The empirical loss function is rewritten as

$$K_n(w) = \frac{1}{2}\|I(u)^{\frac{1}{2}}\Big(u-I(u)^{-1}\frac{\xi_n(u)}{\sqrt{n}}\Big)\|^2$$
$$-\frac{1}{2n}(\xi_n(u)\cdot I(u)^{-1}\xi_n(u)).$$

In the neighborhood of $u=0$, we obtain

$$K_n(w) \cong \frac{1}{2}\|I(0)^{\frac{1}{2}}\Big(u-I(0)^{-1}\frac{\xi_n(0)}{\sqrt{n}}\Big)\|^2$$
$$-\frac{1}{2n}(\xi_n(0)\cdot I(0)^{-1}\xi_n(0)).$$

For an arbitrary function $F(\ )$,

$$\int F(\sqrt{n}\,u)dw$$

$$= \int F(\sqrt{n}\,u) \prod_{j=1}^{g} \delta\Big(u - \prod_{k=d_{j-1}+1}^{d_j} w_k\Big) dw\,du$$

$$= \int F(u) \prod_{j=1}^{g} \delta\Big(\frac{u}{\sqrt{n}} - \prod_{k=d_{j-1}+1}^{d_j} w_k\Big) dw\,\frac{du}{n^{g/2}}$$

$$= \frac{c_3(\log n)^{m-1}}{n^{g/2}} \int F(u)du.$$

On the other hand,

$$
\begin{aligned}
nK_n(w) &= \frac{1}{2}\|I(0)^{1/2}\Big(\sqrt{n}\,u - I(0)^{-1}\xi_n(0)\Big)\|^2 \\
&\quad - \frac{1}{2}(\xi_n(0) \cdot I(0)^{-1}\xi_n(0)) \\
&\equiv \hat{K}_n(\sqrt{n}\,u).
\end{aligned}
$$

Therefore,

$$
\begin{aligned}
\mathbb{E}_w[K_n(w)] &= \frac{\int K_n(w)\exp(-n\beta K_n(w))\varphi(w)dw}{\int \exp(-n\beta K_n(w))\varphi(w)dw} \\
&= \frac{1}{n}\frac{\int \hat{K}_n(\sqrt{n}\,u)\exp(-\beta\hat{K}_n(\sqrt{n}\,u))\varphi(w)dw}{\int \exp(-\beta\hat{K}_n(\sqrt{n}\,u))\varphi(w)dw} \\
&= \frac{1}{n}\frac{\int \hat{K}_n(u)\exp(-\beta\hat{K}_n(u))du}{\int \exp(-\beta\hat{K}_n(u))du} \\
&= \frac{1}{2n}\frac{\int \|I(0)^{\frac{1}{2}}(u - \xi_n^*)\|^2\exp(-\beta\hat{K}_n(u))du}{\int \exp(-\beta\hat{K}_n(u))du} \\
&\quad - \frac{1}{2n}(\xi_n(0) \cdot I(0)^{-1}\xi_n(0)),
\end{aligned}
$$

where the notation

$$\xi_n^* = I(0)^{-1}\xi_n(0)$$

is used. Finally, by the integral formlula

$$\frac{\int \|I(0)^{1/2}u\|^2\exp(-\frac{\beta}{2}\|I(0)^{1/2}u\|^2)du}{\int \exp(-\frac{\beta}{2}\|I(0)^{1/2}u\|^2)du} = \frac{g}{\beta}$$

and by Lemma.4, we have

$$\mathbb{E}[\mathbb{E}_w[K_n(w)]] = \frac{g}{2\beta n} - \frac{g}{2n} + o(\frac{1}{n}),$$

13

|        | regular   | quasi-regular | singular                         |
|--------|-----------|---------------|----------------------------------|
| RCLT   | $d/2$     | $g/2$         | $\lambda$                        |
| SF     | $d/2$     | $g/2$         | $\nu$                            |
| $G_n$  | $d/(2n)$  | $g/(2n)$      | $((\lambda-\nu)/\beta+\nu)/n$    |
| $T_n$  | $-d/(2n)$ | $-g/(2n)$     | $((\lambda-\nu)/\beta-\nu)/n$    |

Table 1: Regular, Quasi-Regular, and Singular

Then, because $\lambda = \frac{g}{2}$ holds from Lemma1, we obtain the Theorem. (Q.E.D.)

**Example.4** By the main theorem of this paper, the real log canonical threshold and the singular fluctuation of Example.1 are $\lambda = \nu = 1$. Also those of Example.3 are $\lambda = \nu = 2$.

# 5 Discusion

Let us discuss the result of this paper from the two different points of view. Firstly, we study the theoretical aspect and then the practical aspect.

## 5.1 Theoretical point of view

In the present paper, we introduced a new concept, a quasi-regular case. A quasi-regular case is not a regular case, but it has the same property as the regular case. Table.1 shows comparison of the real log canonical threshold (RLCT), singular fluctuation (SF), the generalization error $G_n$, and the training error $T_n$.

Even for the general singular cases, real log canonical thresholds have been clarified in several cases. However, this paper is the first case in which the singular fluctuation was clarified. In general singular cases, it is conjectured that the real log canonical threshold is not equal to the singular fluctuation. To clarify such conjecture is the future study.

## 5.2 Practical point of view

In applications, even if both birational invariants are unknown, the generalization error can be estimated from the training error and the functional variance [13] because

$$\mathbb{E}[G_n] = \mathbb{E}[T_n] + \frac{\beta}{n}\mathbb{E}[V_n] + o(\frac{1}{n}),$$

which is asymptotically equivalent to Bayes cross validation [14].

However, in Bayes estimation, the method how to approximate the posterior distribution using Markov chain Monte Carlo (MCMC) method is an important issue. There are a lot of parameters which determine the MCMC process, for example, times of burn-in, times of sufficiently updates, and so on. If we know the concrete values of birational invariants, then we can evaluate how accurate

the MCMC process is [9]. Therefore, the quasi-regular cases are appropriate for evaluating MCMC process. It is the future study to evaluate MCMC process using the quasi-regular cases.

# 6    Conclusion

In the present paper, a new concept, a quasi-regular case, was firstly proposed, and its theoretical foundation was constructed. A quasi-regular case is not a regular case but a singular case, whereas it has the same property as a regular case. In a quasi-regular case, it was proved that the real log canonical threshold is equal to the singular fluctuation. This is the first case in which nontrivial value of singular fluctuation is clarified.

### Acknowledgement

# References

[1] M.Aoyagi,S.Watanabe,"Resolution of singularities and generalization error with Bayesian estimation for layered neural network," Vol.J88-D-II, No.10, pp.2112-2124, 2005.

[2] M.Aoyagi, S.Watanabe,"Stochastic complexities of reduced rank regression in Bayesian estimation," Neural Networks, Vol.18,No.7, pp.924-933, 2005.

[3] M.F.Atiyah,"Resolution of singularities and division of distributions," Comm. Pure Appl. Math., Vol.13,pp.145-150,1970.

[4] K. Hagiwara, "On the Problem in Model Selection of Neural Network Regression in Overrealizable Scenario," Neural Comput., Vol.14,Vol.8, pp.1979 - 2002, 2002.

[5] J.A.Hartigan,"A failure of likelihood asymptotics for normal mixture," Proc. of Barkeley Conf. in honor of Jerzy Neyman and Jack Keifer, Vol.2, pp.807-810,1985.

[6] T. Hayasaka, M. Kitahara, and S. Usui, "On the Asymptotic Distribution of the Least-Squares Estimators in Unidentifiable Models," Neural Comput., Vol.16 ,No.1, pp.99 - 114, 2004.

[7] H. Hironaka, "Resolution of singularities of an algebraic variety over a field of characteristic zero," Ann. of Math., Vol.79, 109-326,1964.

[8] M. Kashiwara, "B-functions and holonomic systems," Inventions Math., 38, 33-53.1976.

[9] K. Nagata, S. Watanabe, "Asymptotic Behavior of Exchange Ratio in Exchange Monte Carlo Method," International Journal of Neural Networks, Vol. 21, No. 7, pp. 980-988, 2008.

[10] S. Watanabe, "Generalized Bayesian framework for neural networks with singular Fisher information matrices," Proc. of International Symposium on Nonlinear Theory and Its applications, (Las Vegas), pp.207-210, 1995.

[11] S. Watanabe, "Algebraic Analysis for Nonidentifiable Learning Machines," Neural Computation, Vol.13, No.4, pp.899-933, 2001.

[12] S. Watanabe, "Algebraic geometrical methods for hierarchical learning machines," Neural Networks, Vol.14, No.8,pp.1049-1060, 2001.

[13] Sumio Watanabe, "Equations of states in singular statistical estimation", Neural Networks, Vol.23, No.1, pp.20-34, 2010.

[14] Sumio Watanabe, "Asymptotic Equivalence of Bayes Cross Validation and Widely Applicable Information Criterion in Singular Learning Theory," Journal of Machine Learning Research, Vol.11, (DEC), pp.3571-3591, 2010.

[15] S. Watanabe, "Algebraic geometry and statistical learning theory," Cambirdge University Press, 2009.

[16] K.Yamazaki, S.Watanabe,"Singularities in mixture models and upper bounds of stochastic complexity." International Journal of Neural Networks, Vol.16, No.7, pp.1029-1038,2003.