

論文97-34C-7-9

# Fokker-Plank 방정식의 해석을 통한 Langevine 경쟁학습의 동역학 분석

(Analysis of the Fokker-Plank Equation for the Dynamics  
of Langevine Competitive Learning Neural Network)

石 鎮 旭 \* , 趙 成 元 \*

(Jinwuk Seok and Seongwon Cho)

## 요 약

본 논문에서는 Langevine 경쟁학습 신경회로망의 Fokker-Plank방정식에 바탕을 둔 동역학 분석을 논한다. 확률미분방정식적 관점에서 Langevine 경쟁학습의 학습알고리즘은 Langevine 확률미분방정식의 한 종류이며 확률측도를 가진 위상공간  $(\Omega, \mathcal{F}, P)$  위의 확산 방정식을 가지고 있다. 제안한 알고리즘에서 먼저 학습방정식의 분석을 통해 학습방정식의 Fokker-Plank 방정식을 유도하여 Markov 준군위의 극소 사용소를 도입하여 특정 단체(simplex)내의 가중치 벡터가 목적함수의 불록, 혹은 의사 불록(pseudo-convex) 조건하에서 대역 최소점을 찾을 수 있음을 증명한다. 원격탐사 데이터에 대한 평균 인식 실험을 통하여 기존의 경쟁학습 신경회로망보다 제안한 Langevine 경쟁학습 신경회로망의 성능이 우수함과 본 논문에서 제기한 대역 기능이 가능함을 증명하였다.

## Abstract

In this paper, we analyze the dynamics of Langevine competitive learning neural network based on its Fokker-Plank equation. From the viewpoint of the stochastic differential equation (SDE), Langevine competitive learning equation is one of Langevine stochastic differential equation and has the diffusion equation on the topological space  $(\Omega, \mathcal{F}, P)$  with probability measure. We derive the Fokker-Plank equation from the proposed algorithm and prove by introducing a infinitesimal operator for Markov semigroups, that the weight vector in the particular simplex can converge to the globally optimal point under the condition of some convex or pseudo-convex performance measure function. Experimental results for pattern recognition of the remote sensing data indicate the superiority of Langevine competitive learning neural network in comparison to the conventional competitive learning neural network.

## I. 서 론

현재까지 신경회로망에 대한 연구는 대부분 특정 적용대상에 대한 직관적인 신경회로망 모델 개발에 중점이 맞추어져 이루어져 왔다. 따라서 범용적으로 쓰이

는 몇몇 모델(오류 역전과 모델이나 자기 조직화 형상지도 모델 혹은 흡혈드, 유전 알고리즘 등)을 제외하고 신경회로망 모델들의 안정성이나 수렴성 혹은 동역학적 연구는 거의 이루어지지 않은 형편이다. 그러므로 각 신경회로망 모델에 대한 동역학(Dynamics)을 파악하여 각 알고리즘이 어떠한 경우에 안정적인 동작을 하는가에 대한 정보와 각 신경회로망 모델들이 어떠한 경우 국소 수렴성을 피할 수 있는 가를 파악하는 것은 공학적으로 매우 중요한 문제라 하겠다.

\* 正會員, 弘益大學校 電子·電氣 工學部

(School of Electronic and Electrical Engineering,  
Hong Ik University)

接受日字: 1996年7月22日, 수정완료일: 1997年7月4日

본 연구에서는 기존의 경쟁학습 신경회로망보다 실 험적으로 학습성능이 우수한 것으로 알려지고 하드웨어 구현에 편리한 Langevine 경쟁학습 신경회로망의 확산성과 대역 최소화 성질을 분석한다<sup>[8]</sup>. 확률 미분 방정식(SDE)의 관점에서 볼 때, Langevine 경쟁학습 신경회로망은 두 가지 방법으로 분석 가능하다. 첫 번째 방법은 학습 방정식의 확률적 요소를 추출해 내어 학습 방정식의 Fokker-Plank 방정식을 유도, 분석하는 방법이며<sup>[3]</sup>, 또 하나의 방법은 학습 방정식을 미분 다양체위의 상 미분 방정식의 해석 기법들을 사용하여 분석하는 방법이다.

첫 번째 방법은 Transition 확률이 준군(Semi-group)을 이루는 특징을 사용, 시간과 위치에 관한 2차 편미분 방정식인 Fokker-Plank 방정식을 유도, 선형 준군동형사상(Linear Semigroup Operator)에 대한 해석을 통해 Weight Vector의 확산성, 안정성, 수렴성 등의 동역학적 특징을 분석하는 방법이며 두 번째 방법은 Langevine 경쟁학습 신경회로망 학습 방정식에 확률적 요소(이진 강화함수)가 있음을 통해 학습 방정식을 확률 미분 방정식으로 놓고 일반적인 상 미분 방정식의 해석 방법을 사용하여 동역학을 파악하는 방법이 있다<sup>[2][4]</sup>.

본 논문에서는 먼저 Langevine 경쟁학습 신경회로망의 학습 알고리즘이 기존의 경쟁학습 알고리즘에 Gaussian 분포를 가진 Brownian Motion Process 가 결합된 것임을 밝히고 이를 통해 Langevine 경쟁학습 신경회로망의 학습 알고리즘의 Fokker-Plank 방정식을 유도한다. 유도된 Fokker-Plank 방정식은 Markov 준군위에 정의된 2차 편 미분 방정식이므로 선형 준군동형사상을 도입하여 Transition 확률의 축소(Contraction)가 어떤 조건 위에 성립하는지를 밝혀 대역 최소성의 조건을 밝힌다. 마지막으로 세 가지 종류의 원격탐사 데이터에 대한 폐턴인식실험을 통해 단위 단체(Simplex)내에 기존의 확률적 최적화 알고리즘으로 구성된 경쟁학습 (Simple Competitive Learning : SCL) 신경회로망과 성능 비교를 통해 Langevine 경쟁학습 신경회로망의 우수성을 확인하였다.

## II. Langevine 경쟁학습 신경회로망 알고리즘

Langevine 경쟁학습 신경회로망의 알고리즘은 다

음과 같다<sup>[1]</sup>.

**step 0** : Weight vector 초기화

**step 1** : Input vector  $x^t$ 에 대하여 다음 조건을 만족하는 Weight vector index  $r$ 을 선택

$$r = \arg \min_{w_i(t)} d(v(t), w_i(t))$$

**step 2** : Input vector  $x^t$ 과 Weight vector index  $r$ 에 대하여 Directional Derivation  $k_t$  계산

$$F(t) := \text{Final Epoch} - t 사이의 임의의 값을 취하는 확률 변수$$

$$P(t) := 0에서 t 사이의 임의의 값을 취하는 확률 변수$$

$$\delta(t) = 1 \quad \text{with probability } \frac{1}{2} P_t, F(t) > P(t)$$

$$\pi(\lambda, \delta, t, v_t, w_t) = \begin{cases} \delta(t) = 1 & -\lambda \delta(t) \operatorname{sgn}(v_t - w_r(t)) \\ \delta(t) = 0 & 0 \end{cases}$$

$$k_t = \epsilon_L \nabla J_t + \pi(\lambda, \delta, t, v_t, w_t) = -\epsilon_L (v_t - w_r(t)) - \lambda \delta(t) \operatorname{sgn}(v_t - w_r(t))$$

**step 3** : Weight vector  $w'$  을 다음 식에 의해 계산

$$w_r(t+1) = w_r(t) + \epsilon_L (v_t - w_r(t)) + \delta \rho(t) \operatorname{sgn}(v_t - w_r(t))$$

**step 4** :  $t$ 를  $t+1$ 로 변환

**step 5** : 한계 Epoch (Final Epoch)까지 진행되었으면 끝. 그렇지 않으면 Step 1.

여기에서  $w'$ 는 시간  $t$ 일 때  $t$ 번째 Weight vector를 의미하며  $\epsilon_L$ 은 시간에 대하여 불변인 일정직등 이득,  $v_t$ 는 시간  $t$ 일 때 입력 Vector, 그리고  $\pi(\lambda, \delta, t, v_t, w')$ 는 이진 강화 함수 항으로서 다음과 같이 구성된다.

$$\pi(\lambda, \delta, t, v_t, w') = -\lambda \delta(t) \operatorname{sgn}(v_t - w') \quad (1)$$

식 (1)에서  $\lambda$ 는 이진 강화함수의 이득이며,  $\delta(t)$ 는 1 혹은 0을 반환하는 Random Variable로서 본 논문에서는 Epoch에 종속인  $P_t$ 의 확률로 1을 반환하며,  $\operatorname{sgn}(v_t - w')$ 은 다음을 만족하는 함수이다.

$$\operatorname{sgn}(v_t - w') = \begin{cases} 1 & v_t - w' \geq 0 \\ -1 & v_t - w' < 0 \end{cases} \quad (2)$$

이때 제안한 경쟁학습 신경회로망의 성능지표  $J_t$ 는 다음과 같다<sup>[5]</sup>.

$$J_t = \frac{1}{2} \int_{B_t} \sum_i \|v_i - w'_i\|^2 dP(v) \quad (3)$$

한편 시간  $t$ 에서의 성능지표  $J_t$ 에 대한 Directional Derivation  $k_t$ 를  $J_t$ 의 Gradient로 놓으면

$$k_t = -\nabla J_t = (v_t - w'_t)$$

따라서 Directional Derivation  $h_i$ 의 Theorem Descent Direction  $h^D_i$ 를 Theorem 1과 같이 놓으면  $h^D_i$ 는  $J_i$ 를 최소화 하는 방향으로  $w^r_i$ 를 보낸다.

Theorem 1 : Descent Direction  $h^D_i$ 가 다음과 같이 정의되면

$$h^D_i = \varepsilon_L h_i + \lambda \delta(t) \operatorname{sgn}(v_i - w^r_i) \quad \varepsilon_L > 0$$

학습 방정식  $w^r_{t+1} = w^r_t + h^D_i$ 는  $J_i$ 를 최소화 하는 방향으로  $w^r_i$ 를 보낸다.

proof :

이진 강화 함수와  $J_i$ 의 Weight vector  $w^r_i$ 에 대한 Gradient  $\nabla J_i$ 와  $h^D_i$ 의 스칼라 적이 0보다 작다는 것과 위 명제는 동치이다. 따라서 이진 강화함수와  $\nabla J_i$ 의 스칼라 적을 살펴보면 이진 강화함수의 정의에 의해 이진 강화함수 n차원 각 Component들은 크기가 1 혹은 0이며 부호가  $\nabla J_i$ 의 대응되는 Component와 반대이므로

$$-n \leq \langle \nabla J_i, \lambda \delta(t) \operatorname{sgn}(v_i - w^r_i) \rangle \leq 0$$

그러므로

$$\langle \nabla J_i, \varepsilon_L h_i + \lambda \delta(t) \operatorname{sgn}(v_i - w^r_i) \rangle = -\varepsilon_L \|\nabla J_i\|^2 + \langle \nabla J_i, \lambda \delta(t) \operatorname{sgn}(v_i - w^r_i) \rangle < 0$$

■

### III. Gaussian Process로서의 이진 강화함수

앞 절 Theorem 1에서 살펴보았듯이 이진 강화함수는 기존의 경쟁학습이 Weight vector  $w^r_i$ 를  $\nabla J_i$ 의 반대 방향으로 개선하는데 비하여 이진 강화함수가 결합된 경우 개선거리가 보다 크고 방향이  $\nabla J_i$  방향에서 약간 틀어진 방향으로 보내게 됨을 알 수 있다. 이진 강화함수는 또한 무작위로 발생하게 되므로 확률로서 학습 방정식의 동역학을 분석할 수 밖에 없는데 이를 위해 본 절에서는 먼저 이진 강화함수가 어떤 프로세스인지를 알아보도록 한다.

Assumption 1 :

i. 입력벡터  $v_i$ 는 최적 Weight vector  $\widehat{w}^r$ 의 Neighborhood 주위에 충분히 조밀하게 분포하고 있다.  
i.e.

$$\forall v_i \in R^n \& \rho_i > 0 \text{ s.t. } v_i \in B^o(\widehat{w}^r, \rho_i)$$

$$= \{ v_i \in R^n | d(v_i, \widehat{w}^r) < \rho_i \} \text{ as } t \rightarrow \infty$$

ii. 입력벡터  $v_i$ 는 Gaussian 분포를 따른다.

iii. 입력벡터  $v_i$ 는 variance  $\|D\| > 0$ 를 가지면서 Identically, Independent Distribution이다.

Assumption 2의 첫 번째 가정은 논의의 대상을 충분히 많은 입력벡터  $v_i$ 에 대하여 제한한 것이고 두 번째 가정은 매우 많은 데이터의 경우 일반적으로 나타나는 분포가 Gaussian 분포이기 때문이다<sup>[6]</sup>. 세 번째 가정은 입력 데이터의 분포가 White Noise에 의해 교란되어 나타나는 것임을 의미한다. 이것은 첫 번째와 두 번째 가정이 성립할 때 일반적으로 나타나는 성질이기도 하다.

증명에 들어가기 앞서 제안한 알고리즘이 가지는 일반적인 성질에 따르는 가정을 두자.

Assumption 2 :

i.  $v_i - w^r_i \geq 0$ 일 확률은  $t \rightarrow \infty$ 에서 다음과 같다.

$$P\{v_i - w^r_i \geq 0\} = 0.5$$

ii.  $t \rightarrow \infty$ 에서 최적 Weight vector  $w^*$  주변에는 평균적으로 강화함수의 영향이 없다.

$$\text{i.e. } \lim_{M \rightarrow \infty} E \sum_{i=0}^M \eta(\lambda, \delta, t, x, w^r_i) < \infty$$

Assumption 1과 Assumption 2에서 다음의 정리 2가 유도된다.

Theorem 2 :

이진 강화함수항  $\eta(\lambda, \delta, t, v_i, w^r_i) = \lambda \delta(t) \operatorname{sgn}(v_i - w^r_i)$ 는 최적 Weight vector  $\widehat{w}^r$ 의 Neighborhood 주위에서 Independent Incremental Gaussian Process이다.

proof :

이진 강화함수항

$\eta(\lambda, \delta, t, v_i, w^r_i) = \lambda \delta(t) \operatorname{sgn}(v_i - w^r_i)$ 에 의해 생성되는 확률 프로세스를  $X_i$ 라 하고 초기 시간  $t_0$  이후, 시간  $t$ 에서 이진 강화 함수항에 의한 변위를  $BR(t)$ 라 놓으면  $BR(t)$ 는 다음과 같다.

$$BR(t) = \lambda(X_{t_0+1} + X_{t_0+2} + \dots + X_{t-1} + X_t)$$

임의  $k$ 번째  $\operatorname{sgn}(v_{t_0+k} - w^r_{t_0+k})$ 의 값이 1 혹은 -1을

가질 확률은 Assumption 2에서 각각 다음과 같다.

$$\begin{aligned} & \operatorname{sgn}(v_{t_0+k} - w^r_{t_0+k}) \\ &= \begin{cases} 1 & \text{with } P(v_{t_0+k} - w^r_{t_0+k} \geq 0 | \delta(t_0+k) = 1) = 0.5 \\ -1 & \text{with } P(v_{t_0+k} - w^r_{t_0+k} \geq 0 | \delta(t_0+k) = 1) = 0.5 \end{cases} \end{aligned}$$

위 식에서  $\delta(t_0+k) = 1$  일 확률을  $P_\delta$ 라 하면  $\delta(t_0+k) = 0$  일 확률은  $1 - P_\delta$ 이다. 따라서

$$\begin{aligned} P(v_{t_0+k} - w^r_{t_0+k} \geq 0 | \delta(t_0+k) = 1) &= \frac{P(v_{t_0+k} - w^r_{t_0+k} \geq 0, \delta(t_0+k) = 1)}{P(\delta(t_0+k) = 1)} \\ &= \frac{1}{2} \end{aligned}$$

$$P(v_{t_0+k} - w^r_{t_0+k} \geq 0, \delta(t_0+k) = 1) = \frac{1}{2} P(\delta(t_0+k) = 1) = \frac{1}{2} P_\delta$$

마찬가지로

$$P(v_{t_0+k} - w^r_{t_0+k} < 0, \delta(t_0+k) = 1) = \frac{1}{2} P(\delta(t_0+k) = 1) = \frac{1}{2} P_\delta$$

따라서  $X_{t_0+k}$ 은 각각 다음의 확률을 가진다.

$$X_{t_0+k} = \begin{cases} 1 & \text{with probability } \frac{1}{2} p_\delta \\ -1 & \text{with probability } \frac{1}{2} p_\delta \\ 0 & \text{with probability } 1 - p_\delta \end{cases}$$

Assumption 1의 iii과 Assumption 2의 i, 그리고  $X_{t_0+k}$ 의 확률에 의하여,  $X_{t_0+k}$ 는  $\forall k \in \mathbb{Z}^+$ 에 대하여 Independent Process이며 따라서  $BR(t)$ 는 Independent Increment 성질을 만족한다. ■

이진 강화 함수형에 의한 변위  $BR(t)$ 가 Gaussian Process임을 증명하기 위하여 먼저  $BR(t)$ 의 평균을 구하면  $X_{t_0+k}$ 가 1을 가질 확률과 -1을 가질 확률이 같으므로  $BR(t)$ 의 평균은 0이고 다음과 같이 표시한다.

$$EBR(t) = 0$$

$BR(t)$ 의 평균은 0이므로 각  $k$ 에 대한  $X_{t_0+k}$ 의 Variance는

$$VAR(X_{t_0+k}) = E X_{t_0+k}^2 - p_\delta$$

따라서  $BR(t)$ 의 Variance는 다음과 같다.

$$\begin{aligned} VAR(BR(t)) &= EBR^2(t) = \lambda^2(t - t_0)p_\delta \\ VAR(BR(t))|_{t_0=0} &= \lambda^2 t p_\delta \end{aligned}$$

한편 Assumption 1의 ii에서 입력  $v_{t_0+k}$ 는  $v_{t_0+k} \in B^o(w^r_{t_0+k})$ 이며 평균  $w^r$ 인 Gaussian 분포를 가지므

로

$$P(v_{t_0+k} - w^r_{t_0+k} \geq 0 | \delta(t_0+k) = 1) = \int_{v_{t_0+k}-w^r_{t_0+k}}^{\infty} \frac{1}{\sqrt{2\pi\sigma_{p_\delta}}} e^{-\frac{(v_{t_0+k}-w^r_{t_0+k})^2}{2\sigma_{p_\delta}^2}} dv_{t_0+k}$$

따라서  $P(v_{t_0+k} - w^r_{t_0+k} \geq 0 | \delta(t_0+k) = 1)$ 의 확률밀도 함수

$$\begin{aligned} & f(v_{t_0+k} - w^r_{t_0+k} \geq 0 | \delta(t_0+k) = 1) \\ &= d_{v_{t_0+k}} P(v_{t_0+k} - w^r_{t_0+k} \geq 0 | \delta(t_0+k) = 1) \end{aligned}$$

는 Gaussian 분포이며, 각  $X_{t_0+k}$ 의 Gaussian 분포를 따르므로 Chain Rule에 의해

$$\begin{aligned} J(BR(t) = B | B_D) &= J^{-1}\left(\frac{\partial BR(t)}{\partial X_{t_0+k}}\right) J^{-1}\left(\frac{\partial X_{t_0+k}}{\partial v_{t_0+k}}\right) \\ &\cdot f(v_{t_0+k} - w^r_{t_0+k} \geq 0 | \delta(t_0+k) = 1) \end{aligned}$$

$BR(t)$ 는 평균 0, Variance  $\lambda^2 t p_\delta$ 를 따르는 Gaussian 분포를 가진다.

그러므로  $BR(t)$ 는 Gaussian Process이다. ■

Corollary :  $BR(t)$ 의 Co-variance는  $s \leq t$ 에 대하여  $COV(BR(s), BR(t)) = \lambda^2 p_\delta \delta(t-s)$ 이다.

Proof :

$$\begin{aligned} COV(BR(s), BR(t)) &= COV(BR(s), BR(s) + BR(t) - BR(s)) \\ &= COV(BR(s), BR(s)) + COV(BR(s), BR(t) - BR(s)) \\ &= COV(BR(s), BR(s)) = \lambda^2 p_\delta \delta(t-s) \end{aligned}$$

#### IV. Langvine 경쟁학습 신경회로망의 Fokker-Plank 방정식 유도

Theorem 2에서 이진 강화함수 항은 Gaussian Process임이 밝혀졌으므로, 이진 강화함수와 일정적응 이득을 가진 경쟁학습 신경회로망의 학습 방정식은 다음과 같이 일반적인 Langevine 방정식의 형태를 가지게 된다.

$$\partial w^r_t = -\nabla J(w^r_t)dt + \eta(t, \cdot)dt \quad (4)$$

Theorem 2에서 이진 강화함수 항  $\eta(t, \cdot)$ 은 "Gaussian Process"으로 다음의 성질을 만족한다.

$$\begin{aligned} E(\eta(t, \cdot)) &= \langle \eta(t, \cdot) \rangle = 0 \\ E(\eta(t, \cdot)\eta(t, \cdot)) &= \langle \eta(t, \cdot)\eta(s, \cdot) \rangle = \lambda^2 p_\delta \delta(t-s) \\ E\left(\prod_{i=1}^n \eta(t_i, \cdot)\right) &= \langle \prod_{i=1}^n \eta(t_i, \cdot) \rangle = \begin{cases} 0 & : l \text{ is odd} \\ \sum \prod \lambda^2 p_\delta \delta(t-s) & : l \text{ is even} \end{cases} \end{aligned}$$

시간  $t$ 에서  $t$ 번째 Weight vector  $w'_t$ 이 열린구간  $(w'_t, w'_t + dw'_t)$  사이에 있을 확률에 대한 확률밀도 함수를  $\rho(w'_t, t)$ 라 놓고  $\rho(w'_t, t)$ 의 시간에 대한 변화율을 고려하면,

$$\frac{\partial \rho(w'_t, t)}{\partial t} = -\frac{\partial}{\partial w'_t} \left[ \frac{\partial w'_t}{\partial t} \rho(w'_t, t) \right] \quad (5)$$

식 (1)을 식 (2)에 대입하면

$$\frac{\partial \rho(w'_t, t)}{\partial t} = -\frac{\partial}{\partial w'_t} [(-\nabla J(w') + g(t, \cdot)) \rho(w'_t, t)] \quad (6)$$

여기서 변위에 대한 Foward 선형 미분 연산자  $L_t$ 를 다음과 같이 정의한다.

**Definition 1:** 변위에 대한 Foward 선형 미분 연산자  $L_t$ 는 다음과 같이 정의된다.

$$L_t \rho(x, t) \equiv \frac{\partial}{\partial x} [A(x) \rho(x, t)] = \left[ \frac{\partial A(x)}{\partial x} + A(x) \frac{\partial}{\partial x} \right] \rho(x, t)$$

Definition 1에 의해 정의된 연산자  $L_t$ 를 식 (6)에 도입하면

$$\begin{aligned} L_t \rho(w'_t, t) &= \frac{\partial}{\partial w'_t} [\nabla J(w') \rho(w'_t, t)] \\ &= \left[ \frac{\partial \nabla J(w')}{\partial w'_t} + \nabla J(w') \frac{\partial}{\partial w'_t} \right] \rho(w'_t, t) \end{aligned} \quad (7)$$

식 (7)을 식 (6)에 대입하면

$$\frac{\partial \rho(w'_t, t)}{\partial t} = L_t \rho(w'_t, t) - \frac{\partial \rho(w'_t, t)}{\partial w'_t} g(t, \cdot) \quad (8)$$

편미분 방정식 (8)을 시간에 대하여 풀기 위해 식 (8)의 해의 한 후보를 다음과 같이 놓는다.

$$\rho(w'_t, t) = e^{tL} \xi(w'_t, t) \quad (9)$$

식 (9)에서  $\xi(w'_t, t)$ 는 시간  $t$ 에 종속인 함수이므로  $w'_t$ 에 대한 미분 값은 0이다. 식 (9)을  $t$ 에 대하여 편미분하면

$$\begin{aligned} \frac{\partial \xi(w'_t, t)}{\partial t} &= -L_t(w') \rho(w'_t, t) + e^{-tL} \frac{\partial \rho(w'_t, t)}{\partial t} \\ &= e^{-tL} [-L_t(w') \rho(w'_t, t) + L_t(w') \rho(w'_t, t) - g(t, \cdot) \frac{\partial \rho(w'_t, t)}{\partial w'_t}] \\ &= -g(t, \cdot) e^{-tL} \left[ 1 - \frac{\partial e^{-tL}}{\partial w'_t} \right] \xi(w'_t, t) + \frac{\partial \xi(w'_t, t)}{\partial w'_t} e^{-tL} \end{aligned} \quad (10)$$

식 (10)에서  $\xi(w'_t, t)$ 의  $t$ 에 대한 미분연산자  $K(w'_t, t)$ 를 다음과 같이 놓는다.

$$K(w'_t, t) = -g(t, \cdot) e^{-tL} \left[ 1 - \frac{\partial e^{-tL}}{\partial w'_t} \right] \xi(w'_t, t) - \frac{\partial \xi(w'_t, t)}{\partial w'_t} e^{-tL}$$

$\xi(w'_t, t)$ 를 적분형 Taylor급수로 놓으면

$$\begin{aligned} \xi(w'_t, t + \Delta t) &= \xi(w'_t, t) + \int_t^{t+\Delta t} dt_1 K(w'_t, t) \xi(w'_t, t) \\ &= [1 + \int_t^{t+\Delta t} dt_1 K_1(w'_t, t) + \frac{1}{2!} \int_t^{t+\Delta t} dt_1 \int_t^{t+\Delta t} dt_2 K_2(w'_t, t) K_2(w'_t, t) \\ &\quad + o(w'_t, t)] \xi(w'_t, t) \end{aligned} \quad (11)$$

$\xi(w'_t, t + \Delta t)$ 의 평균을 식 (11)에서 구하게 되면

$$\begin{aligned} E\xi(w'_t, t + \Delta t) &= [1 + \int_t^{t+\Delta t} dt_1 EK_1(w'_t, t) + \frac{1}{2!} \int_t^{t+\Delta t} dt_1 \int_t^{t+\Delta t} dt_2 EK_2(w'_t, t) \\ &\quad + o(w'_t, t)] E\xi(w'_t, t) \end{aligned} \quad (12)$$

식 (12)에서  $K(w'_t, t)$ 의 2차 모멘트 까지의 평균을 구하면

$$\begin{aligned} E(K(w'_t, t) \xi(w'_t, t)) &= Eg(t, \cdot) e^{-tL} \left( \frac{\partial}{\partial w'_t} \right) e^{tL} E\xi(w'_t, t) = 0 \\ E(K(w'_t, t) K(w'_t, t) \xi(w'_t, t)) &= Eg(t, \cdot) e^{-tL} \left( \frac{\partial}{\partial w'_t} \right) e^{-tL} \left( \frac{\partial}{\partial w'_t} \right) e^{tL} E\xi(w'_t, t) \\ &\quad + e^{-tL} \left( \frac{\partial}{\partial w'_t} \right) e^{-tL} \left( \frac{\partial}{\partial w'_t} \right) e^{tL} E\xi(w'_t, t) \\ &= \lambda^2 p_g \delta(t-s) e^{-2tL} \left( \frac{\partial^2}{\partial^2 w'_t} \right) e^{tL} E\xi(w'_t, t) \end{aligned}$$

위 결과를 통해 Taylor급수의 2차항에 대한 평균을 구하면

$$\begin{aligned} &\frac{1}{2!} \int_t^{t+\Delta t} dt_1 \int_t^{t+\Delta t} dt_2 E(K_1(w'_t, t) K_2(w'_t, t) \xi(w'_t, t)) \\ &= \frac{1}{2!} \int_t^{t+\Delta t} dt_1 \int_t^{t+\Delta t} \lambda^2 p_g \delta(t_1 - t_2) e^{-t_1 L} \left( \frac{\partial}{\partial w'_t} \right) e^{t_1 L} \left( \frac{\partial}{\partial w'_t} \right) e^{-t_2 L} \left( \frac{\partial}{\partial w'_t} \right) e^{t_2 L} E\xi(w'_t, t) \\ &= \frac{1}{2} \lambda^2 p_g \Delta t e^{-tL} \left( \frac{\partial^2}{\partial^2 w'_t} \right) e^{tL} E\xi(w'_t, t) \end{aligned}$$

한편 고차항에서는 앞에서 살펴본 바와 같이 차수  $n$ 이 홀수이면 0, 2n이면  $(\Delta t)^n$ 에 비례한다. 따라서  $\Delta t \rightarrow 0$ 의 극한의 경우 2차항까지만 고려한다면

$$\begin{aligned} \lim_{\Delta t \rightarrow 0} \frac{E\xi(w'_t, t + \Delta t) - E\xi(w'_t, t)}{\Delta t} &= \frac{\partial E\xi(w'_t, t)}{\partial t} \\ &= \frac{1}{2} \lambda^2 p_g e^{-tL} \left( \frac{\partial^2}{\partial^2 w'_t} \right) e^{tL} E\xi(w'_t, t) \end{aligned}$$

$P(w'_t, t)$ 를  $\rho(w'_t, t)$ 의 ensemble에서 취한 평균으로 놓으면

$$P(w'_t, t) = E\rho(w'_t, t) = e^{tL} E\xi(w'_t, t)$$

따라서 다음의 Fokker-Plank방정식을 얻을 수 있다.

$$\begin{aligned} \frac{\partial P(w'_t, t)}{\partial t} &= L_t e^{tL} E\xi(w'_t, t) + e^{tL} \frac{\partial E\xi(w'_t, t)}{\partial t} \\ &= L_t P(w'_t, t) + \frac{1}{2} \lambda^2 p_g \left( \frac{\partial^2}{\partial^2 w'_t} \right) e^{tL} E\xi(w'_t, t) \\ &= \frac{\partial}{\partial w'_t} [(\nabla J(w')) P(w'_t, t)] + \frac{1}{2} \lambda^2 p_g \left( \frac{\partial^2}{\partial^2 w'_t} \right) P(w'_t, t) \end{aligned}$$

## V. Fokker-Plank 방정식을 사용한 Langevine 경쟁학습 신경회로망의 대역 최소화 근사 해석

앞에서 살펴본 바 Assumption 1에서 “최적 Weight vector  $\hat{w}_t$ 의 Neighborhood을 도입함에 따라 실제로 4절까지의 모든 분석은 국소적인 분석에 지나지 않는다. 따라서 4절에 유도된 Fokker-Plank 방정식의 해를 통한 고찰도 실제로는 국소적인 특성만을 분석한 것이 된다. 그러나 경쟁학습 신경회로망에서는 다수의 Weight Vector들에 의해 입력 벡터공간을 다수의 단체(Simplex)들로 분할하는 것으로 하나의 단체를 덮을 수 있는 Neighborhood을 도입하여, Neighborhood내의 대역 최소점을 찾을 수 있는가를 분석하여야 한다. 이를 위해 먼저 Fokker-Plank방정식에 대한 Infinitesimal Operator  $\mathcal{L}_t^*$  도입하고,  $\mathcal{L}_t^*$ 의 Quasi-Norm을 정의하여 Weight vector  $w_t^*$  Transition 확률에 대한 방정식 (Langevine 경쟁학습 신경회로망의 Fokker-Plank방정식)  $\partial_t P(w_t^*, t) = \mathcal{L}_t^* P(w_t^*, t)$ 이 Contraction되어 있는가를 살펴 보아야 한다. 일반적으로 확산 방정식(Diffusion Equation) 혹은 Fokker-Plank방정식의 해(Transition 확률)가 어떤 조건아래 Contraction 되어 있고 Transition 확률이 성능지표  $J(w_t)$ 의 미분동상 사상(Diffeomorphism)이면 Contraction 조건을 만족하는 범위 내에서 대역 최소점을 찾을 수 있음이 알려져 있다<sup>[7]</sup>.

**Definition 2 :** Fokker-Plank방정식에 대한 Infinitesimal Operator  $\mathcal{L}_t^*$ 을 다음과 같이 정의한다.

$$\mathcal{L}_t^* = -\frac{\partial}{\partial w_t} a(w_t^*, v_t) + \frac{1}{2} \sigma(t) \frac{\partial^2}{\partial^2 w_t}$$

$a(w_t^*, v_t) = -\varepsilon \nabla J(w_t)$ ,  $\sigma(t) = \lambda^2 p_s(t)$ ,  $t > 0$ 이며 Fokker-Plank방정식은

$$\frac{\partial P(w_t^*, t)}{\partial t} = \mathcal{L}_t^* P(w_t^*, t) = 0$$

**Definition 3 :** Infinitesimal Operator  $\mathcal{L}_t^*$ 는  $P(\cdot) \in L[0, 1]$ 인 임의의 함수  $f$ 에 대하여 다음과 같이 Quasi-Norm을 정의한다.

$$\mathcal{L}_t = \|\mathcal{L}\| := \sup_{t_0 \leq t \leq t_1} \mathcal{L}_t^* P(w_t^*, t)$$

Infinitesimal Operator  $\mathcal{L}_t^*$ 의 Quasi-Norm  $\mathcal{L}_t$ 는 오직 다음의 Norm 성질만 만족한다.

$$\mathcal{L}_t = 0 \quad \forall P(w_t^*, t) \in L[0, 1] : \sup_{t_0 \leq t \leq t_1} \mathcal{L}_t^* P(w_t^*, t) = 0$$

**Theorem 3 :** Langevine 경쟁학습 신경회로망의 Fokker-Plank방정식  $\partial_t P(w_t^*, t) = \mathcal{L}_t^* P(w_t^*, t)$  은 Infinitesimal Operator  $\mathcal{L}_t^*$ 의 Quasi-Norm의 시간에 대한 적분이 다음을 만족하면  $L[0, 1]$ 에서 Contraction 되어 있다

$$\exists t_0, t \in R^+ \text{ s.t. } \int_{t_0}^t \mathcal{L}_s ds \leq \log(1 - \frac{\varepsilon}{P(w_t^*, t_0)}) \\ \forall \varepsilon \in (0, P(w_t^*, t_0))$$

*Proof :*

Langevine 경쟁학습 신경회로망의 Fokker-Plank방정식은 다음과 같다.

$$\frac{\partial P(w_t^*, t)}{\partial t} = \mathcal{L}_t^* P(w_t^*, t) \leq \mathcal{L}_t P(w_t^*, t) \quad (13)$$

$\mathcal{L}_t$ 를 다음과 같이 정의하고

$$\mathcal{L}_s = \sup_{s-\epsilon < x < s+\epsilon} \mathcal{L}_s^* P(w_t^*, s) \quad \forall P(w_t^*, s) \in L[0, 1]$$

식 (13)의 양변에  $\exp(-\int_{t_0}^t \mathcal{L}_s ds)$ 을 곱한다.

$$\frac{\partial P(w_t^*, t)}{\partial t} \exp(-\int_{t_0}^t \mathcal{L}_s ds) = \mathcal{L}_t P(w_t^*, t) \exp(-\int_{t_0}^t \mathcal{L}_s ds) \leq 0 \\ \frac{\partial}{\partial t} P(w_t^*, t) \exp(-\int_{t_0}^t \mathcal{L}_s ds) \leq 0 \quad (14)$$

식 (14)를 Lebesgue 적분하면

$$\int_{P(w_t^*, t_0)}^{P(w_t^*, t)} \partial P(w_t^*, t) \exp(-\int_{t_0}^t \mathcal{L}_s ds) \leq 0 \\ P(w_t^*, t) \exp(-\int_{t_0}^t \mathcal{L}_s ds) \leq P(w_t^*, t_0) \\ P(w_t^*, t) \leq P(w_t^*, t_0) \exp(\int_{t_0}^t \mathcal{L}_s ds) \quad (15)$$

식 (15)의 양변에 초기 Transient 확률  $P(w_t^*, t_0)$ 을 빼고 Norm을 취하면 다음의 관계식이 성립한다.

$$P(w_t^*, t) - P(w_t^*, t_0) \leq P(w_t^*, t_0) \exp(\int_{t_0}^t \mathcal{L}_s ds) - P(w_t^*, t_0) \\ \|P(w_t^*, t) - P(w_t^*, t_0)\| \leq \|P(w_t^*, t_0)\| \exp(\int_{t_0}^t \mathcal{L}_s ds) - 1 \\ \leq \|P(w_t^*, t_0)\| \cdot \|1 - \exp(\int_{t_0}^t \mathcal{L}_s ds)\|$$

가정에서  $\int_{t_0}^t \mathcal{L}_s ds \leq \log(1 - \frac{\varepsilon}{P(w_t^*, t_0)})$  이므로

$$\|P(w_t^*, t) - P(w_t^*, t_0)\| \leq \|P(w_t^*, t_0)\| \cdot \|1 - \exp(\int_{t_0}^t \mathcal{L}_s ds)\| \\ \leq \|P(w_t^*, t_0)\| \cdot \|1 - (1 - \frac{\varepsilon}{P(w_t^*, t_0)})\| \\ \leq \varepsilon \quad (16)$$

가정에서  $\varepsilon$ 은  $\varepsilon \in (0, P(w_{t_0}^r, t_0))$ 를 만족하는 임의의 양수이므로 Langevine 경쟁학습 신경회로망의 Fokker-Plank방정식  $\partial_t P(w_t^r, t) = \mathcal{L}_t P(w_t^r, t)$ 은  $L[0, 1]$ 에서 Contraction되어 있다. ■

Lemma 1 :

$$\forall \xi > 0, \exists t_\xi > t_0 \text{ s.t. } \int_{t_0}^t \frac{(w_\tau^r - w_{t_0}^r)}{\|w_\tau^r - w_{t_0}^r\|} \eta(\tau, \cdot) P(w_\tau^r, t) d\tau < \xi \quad (17)$$

proof :

Transition 확률  $P(w_t^r, t)$ 은 1 보다 작으므로 다음의 성립한다.

$$\int_{t_0}^t \frac{(w_\tau^r - w_{t_0}^r)}{\|w_\tau^r - w_{t_0}^r\|} \eta(\tau, \cdot) P(w_\tau^r, t) d\tau \leq \int_{t_0}^t \frac{(w_\tau^r - w_{t_0}^r)}{\|w_\tau^r - w_{t_0}^r\|} \eta(\tau, \cdot) d\tau$$

논의를 간편히 하기 위해  $y_\tau = \frac{(w_\tau^r - w_{t_0}^r)}{\|w_\tau^r - w_{t_0}^r\|}$ 이라 놓으면,  $y_\tau$ 의 상한은 1이다. 따라서

$$\int_{t_0}^t \frac{(w_\tau^r - w_{t_0}^r)}{\|w_\tau^r - w_{t_0}^r\|} \eta(\tau, \cdot) d\tau = \int_{t_0}^t y_\tau \eta(\tau, \cdot) d\tau \quad (18)$$

Assumption 1과 2, 그리고 Theorem 2에서

$$\eta(\tau, \cdot) = \begin{cases} \lambda & v_\tau - w_\tau^r \geq 0 \& \delta(\tau) = 1 \\ -\lambda & v_\tau - w_\tau^r < 0 \& \delta(\tau) = 1 \end{cases}$$

혹은

$$\eta(\tau, \cdot) = \begin{cases} \lambda & \text{with Probability } 0.5 E\delta(t) \\ -\lambda & \text{with Probability } 0.5 E\delta(t) \end{cases}$$

이므로

$$\begin{aligned} \int_{t_0}^t y_\tau \eta(\tau, \cdot) d\tau &= \lambda \int_{t_0}^t y_\tau d\tau P(v_\tau - w_\tau^r \geq 0, \delta(\tau) = 1) \\ &\quad - \lambda \int_{t_0}^t y_\tau d\tau P(v_\tau - w_\tau^r < 0, \delta(\tau) = 1) \end{aligned} \quad (19)$$

그런데  $P(\delta(\tau) = 1) \rightarrow 0$  ( $\tau \rightarrow \infty$ )에 따라  $P(\delta(\tau) = 1) \rightarrow 0$  이므로

“Weak Law of Large Numbers”에 의해<sup>[8]</sup>

$$\begin{aligned} \exists t_\xi > t_0 \text{ s.t. } P[y_\tau | y_\tau| \geq \xi] &= \int_{t_0}^t y_\tau d\tau P(v_\tau - w_\tau^r \geq 0, \delta(\tau) = 1) \\ &\quad - \int_{t_0}^t y_\tau d\tau P(v_\tau - w_\tau^r < 0, \delta(\tau) = 1) \\ &= \frac{1}{2} \int_{t_0}^t y_\tau d\tau P(\delta(\tau) = 1) - \frac{1}{2} \int_{t_0}^t y_\tau d\tau P(\delta(\tau) = 1) \\ &= \frac{1}{\lambda} \int_{t_0}^t y_\tau \eta(\tau, \cdot) d\tau < \xi \end{aligned} \quad (20)$$

따라서  $\xi > 0$ 인 임의의  $\xi$ 에 대하여  $\int_{t_0}^t \frac{(w_\tau^r - w_{t_0}^r)}{\|w_\tau^r - w_{t_0}^r\|} \eta(\tau, \cdot) P(w_\tau^r, t) d\tau < \xi$ 을 만족하는  $t_\xi > t_0$ 인  $t_\xi$ 가 존재한다.

다. ■

Assumption 3 : 성능지표  $J(w_t^r)$ 는 Lipschitz 조건을 만족한다 i.e.

$$\exists K > 0 \text{ s.t. } \|J(w_{t_1}^r) - J(w_{t_2}^r)\| \leq K \|w_{t_1}^r - w_{t_2}^r\| \forall J(w_t^r) \in R, \forall w_t^r \in R$$

Lemma 2 :  $R^+$   $\forall \varepsilon \in (0, P(w_{t_0}^r, t_0))$ 에 대하여 성능지표  $J(w_t^r)$ 가 Convex함수일 때 다음을 만족하는  $t > t_0$ 이며  $R^+$ 의 한 원소인  $t$ 가 존재한다.

$$\int_{t_0}^t \mathcal{L}_s ds \leq \log(1 - \frac{\varepsilon}{P(w_{t_0}^r, t_0)}) \quad (21)$$

proof :

위 명제가 거짓이다 즉,  $\forall t > t_0, \int_{t_0}^t \mathcal{L}_s ds > \log(1 - \frac{\varepsilon}{P(w_{t_0}^r, t_0)})$ 이라 가정하자.

Langevine 경쟁학습 신경회로망의 확률미분방정식(SDE)형식에서

$$\begin{aligned} \partial_t w_t^r &= -\varepsilon_L \nabla J(w_t^r) dt + \eta(t, \cdot) dt \\ w_t^r - w_{t_0}^r &= -\int_{t_0}^t \varepsilon_L \nabla J(w_\tau^r) - \eta(\tau, \cdot) d\tau \end{aligned} \quad (22)$$

양변에  $(w_t^r - w_{t_0}^r) P(w_t^r, t)$ 을 곱하면

$$\begin{aligned} \|w_t^r - w_{t_0}^r\|^2 P(w_t^r, t) &= -\int_{t_0}^t [\varepsilon_L (w_\tau^r - w_{t_0}^r) \nabla J(w_\tau^r) \\ &\quad - (w_\tau^r - w_{t_0}^r) \eta(\tau, \cdot)] P(w_\tau^r, \tau) d\tau \end{aligned} \quad (23)$$

Theorem 1과 Convex조건에 의해 Langevine 경쟁학습에 의해 생성되는 수열  $\{w_t^r\}_{t_0}^t$ 에 대하여  $\forall \tau > t_0$

$$J(w_{t_0}^r) - J(w_\tau^r) \geq (w_{t_0}^r - w_\tau^r) \nabla J(w_\tau^r) \circ$$

$$\begin{aligned} \|w_t^r - w_{t_0}^r\|^2 P(w_t^r, t) &= \int_{t_0}^t [\varepsilon_L (w_\tau^r - w_{t_0}^r) \nabla J(w_\tau^r) + (w_\tau^r - w_{t_0}^r) \eta(\tau, \cdot)] P(w_\tau^r, \tau) d\tau \\ &\leq \int_{t_0}^t [\varepsilon_L (J(w_{t_0}^r) - J(w_\tau^r)) + (w_\tau^r - w_{t_0}^r) \eta(\tau, \cdot)] P(w_\tau^r, \tau) d\tau \\ &\leq \int_{t_0}^t [\varepsilon_L \|J(w_{t_0}^r) - J(w_\tau^r)\| + (w_\tau^r - w_{t_0}^r) \eta(\tau, \cdot)] P(w_\tau^r, \tau) d\tau \end{aligned} \quad (24)$$

또한 Lipschitz 조건에서

$$\begin{aligned} \int_{t_0}^t [\varepsilon_L \|J(w_{t_0}^r) - J(w_\tau^r)\| + (w_\tau^r - w_{t_0}^r) \eta(\tau, \cdot)] P(w_\tau^r, \tau) d\tau &\leq \int_{t_0}^t [\varepsilon_L K \|w_\tau^r - w_{t_0}^r\| + (w_\tau^r - w_{t_0}^r) \eta(\tau, \cdot)] P(w_\tau^r, \tau) d\tau \end{aligned} \quad (25)$$

그러므로

$$\begin{aligned} \|w_t^r - w_{t_0}^r\|^2 P(w_t^r, t) &\leq \int_{t_0}^t [\varepsilon_L K \|w_\tau^r - w_{t_0}^r\| \\ &\quad + (w_\tau^r - w_{t_0}^r) \eta(\tau, \cdot)] P(w_\tau^r, \tau) d\tau \end{aligned}$$

따라서

$$\|w_t^r - w_{t_0}^r\| P(w_t^r, t) \leq \int_{t_0}^t \left[ \varepsilon_L K + \frac{(w_t^r - w_{t_0}^r)}{\|w_t^r - w_{t_0}^r\|} \eta(\tau, \cdot) \right] P(w_\tau^r, \tau) d\tau \quad (26)$$

식 (15)에서  $P(w_{t_0}^r, t) \leq P(w_{t_0}^r, t_0) \exp(\int_{t_0}^t \mathcal{L}_s ds)$  이고  
 $P(w_{t_0}^r, t_0) \exp(\int_{t_0}^t \mathcal{L}_s ds) > 0$  이므로

$$\begin{aligned} \|w_t^r - w_{t_0}^r\| P(w_t^r, t) &\leq \int_{t_0}^t \varepsilon_L K P(w_{t_0}^r, t_0) \exp(\int_{t_0}^\tau \mathcal{L}_s ds) d\tau \\ &+ \int_{t_0}^t \frac{(w_\tau^r - w_{t_0}^r)}{\|w_\tau^r - w_{t_0}^r\|} \eta(\tau, \cdot) P(w_\tau^r, \tau) d\tau \quad (27) \\ &\leq \varepsilon_L K P(w_{t_0}^r, t_0) \int_{t_0}^t \exp(\int_{t_0}^\tau \mathcal{L}_s ds) d\tau \\ &+ \int_{t_0}^t \frac{(w_\tau^r - w_{t_0}^r)}{\|w_\tau^r - w_{t_0}^r\|} \eta(\tau, \cdot) P(w_\tau^r, \tau) d\tau \end{aligned}$$

가정에서  $\forall t > t_0, \int_{t_0}^t \mathcal{L}_s ds > \log(1 - \frac{\varepsilon}{P(w_{t_0}^r, t_0)})$  이므로  
 $\int_{t_0}^t \exp(\int_{t_0}^\tau \mathcal{L}_s ds) d\tau$ 의 하한은

$$\int_{t_0}^t \exp(\int_{t_0}^\tau \mathcal{L}_s ds) d\tau > \int_{t_0}^t (1 - \frac{\varepsilon}{P(w_{t_0}^r, t_0)}) d\tau = (1 - \frac{\varepsilon}{P(w_{t_0}^r, t_0)})(t - t_0) \quad (28)$$

이다. Lemma 1에서  $\lim_{t \rightarrow \infty} \int_{t_0}^t \frac{(w_\tau^r - w_{t_0}^r)}{\|w_\tau^r - w_{t_0}^r\|} \eta(\tau, \cdot) P(w_\tau^r, \tau) d\tau \rightarrow 0$  이므로

$$\begin{aligned} \|w_t^r - w_{t_0}^r\| P(w_t^r, t) &> \int_{t_0}^t \frac{(w_\tau^r - w_{t_0}^r)}{\|w_\tau^r - w_{t_0}^r\|} \eta(\tau, \cdot) P(w_\tau^r, \tau) d\tau = H(t_\xi) \\ &\geq \int_{t_0}^{t_\xi + t_0} \frac{(w_\tau^r - w_{t_0}^r)}{\|w_\tau^r - w_{t_0}^r\|} \eta(\tau, \cdot) P(w_\tau^r, \tau) d\tau \quad \forall t_0 > 0 \end{aligned}$$

되는  $t_\xi$ 를 취할 수 있고, 식 (27)은 다음과 같아 쓸 수 있다.

$$\frac{\|w_t^r - w_{t_0}^r\| P(w_t^r, t)}{\varepsilon_L K P(w_{t_0}^r, t_0)} \leq \int_{t_0}^t \exp(\int_{t_0}^\tau \mathcal{L}_s ds) d\tau + \frac{H(t_\xi + t_0)}{\varepsilon_L K P(w_{t_0}^r, t_0)}$$

그런데  $\varepsilon \rightarrow P(w_{t_0}^r, t_0)$ 에 따라  $\inf \int_{t_0}^t \exp(\int_{t_0}^\tau \mathcal{L}_s ds) d\tau \rightarrow 0$   
 $\forall t > t_0, t \in R^+$  이므로

$$\frac{\|w_t^r - w_{t_0}^r\| P(w_t^r, t)}{\varepsilon_L K P(w_{t_0}^r, t_0)} - \frac{H(t_\xi + t_0)}{\varepsilon_L K P(w_{t_0}^r, t_0)} > \int_{t_0}^t \exp(\int_{t_0}^\tau \mathcal{L}_s ds) d\tau \quad (29)$$

되는  $t$ 가 존재한다. 이는 가정에 모순. 따라서 다음을 만족하는  $t > t_0$ 인  $t$ 가 존재한다.

$$\int_{t_0}^t \mathcal{L}_s ds \leq \log(1 - \frac{\varepsilon}{P(w_{t_0}^r, t_0)}). \quad \blacksquare$$

Corollary 2 :  $\forall \varepsilon \in (0, P(w_{t_0}^r, t_0))$ 에 대하여 성능지표  $J(w^r)$ 가 다음과 같은 Pseudo-Convex 함수일 때 식 (21)을 만족하는  $t > t_0$ 이며  $R^+$ 의 한 원소인  $t$ 가 존재한다.

$$\forall \tau > t_0 \quad J(w_{t_0}^r) - J(w_\tau^r) \geq (w_{t_0}^r - w_\tau^r) \nabla J(w_\tau^r) - o(\|w_\tau^r\|) \quad (30)$$

여기서  $o(\|w_\tau^r\|)$ 은  $(w_\tau^r + dw^r)$  범위에서 (30)을 만족할 수 있는  $o(\|w_\tau^r\|)$ 의 절대값의 하한값을 가지며 다음과의 두 특성을 만족한다.

$$\int_{t_0}^t o(\|w_\tau^r\|) P(w_\tau^r, \tau) d\tau \leq 0 \quad o(\|w_\tau^r\|) \rightarrow 0 \text{ as } w_\tau^r \rightarrow w_*^r,$$

$w_*^r$ 는 대역 최소점(Globally Optimal point)이다.

proof :

Lemma 2의 증명에서  $o(\|w_\tau^r\|)$  항을 포함하여 정리하면

$$\begin{aligned} \|w_t^r - w_{t_0}^r\| P(w_t^r, t) &\leq \varepsilon_L K P(w_{t_0}^r, t_0) \int_{t_0}^t \exp(\int_{t_0}^\tau \mathcal{L}_s ds) d\tau \\ &+ \int_{t_0}^t o(\|w_\tau^r\|) P(w_\tau^r, \tau) d\tau \\ &+ \int_{t_0}^t \frac{(w_\tau^r - w_{t_0}^r)}{\|w_\tau^r - w_{t_0}^r\|} \eta(\tau, \cdot) P(w_\tau^r, \tau) d\tau \end{aligned} \quad (31)$$

가정에서  $\int_{t_0}^t o(\|w_\tau^r\|) P(w_\tau^r, \tau) d\tau \leq 0$ 의 최대값은 0이므로 (31)은 (27)과 동일해 진다. 이후의 증명은 Lemma 2의 증명과 동일 ■

Theorem 3, Lemma 2, Corollary 2의 결론에서 성능지표  $J(w^r)$ 가 단위 단체(Simplex)에서 Convex이거나, Corollary 2의 제한조건을 만족하는 Pseudo-Convex 함수의 경우 단위 단체(Simplex)내의 대역 최소점(Globally Optimal point)을 취할 수 있음을 알 수 있다.

## VI. 실험결과

실험에 사용한 데이터는 Flight Line C1(FLC1)으로 불리는 다중분광 지상관측 원격탐사 데이터(Multispectral Earth Observational Remote Sensing Data)로서 미국 인디애나주 Tippecanoe Country 남부지역의 농작물 재배지역을 촬영한 것이다. 이 데이터들은 256 Gray Level로 표현되고 4개의 클래스가 4개의 주요 농산물을 대표하도록 선택된 하나의 데이터와 8개의 클래스가 8개의 주요 농산물을 대표하도록 선택된 데이터 두 가지 종류가 있다. 각 클래스는 학습을 위한 200개의 데이터 벡터들과 테스트를 위한 375

개의 테스트 벡터를 가지며 하나의 데이터 혹은 테스트 벡터는 8차원으로 구성되어 있으며 추정된 확률 분포와 실제 확률분포가 근사하여 분류 결과가 비교적 우수한 특징을 지닌다. 실험에 사용한 또하나의 데이터 집합은 미국 Colorado주의 산악 지역의 지형을 활용한 것으로서 모두 13개의 서로 다른 지형을 나타내며 256 Gray Level로 표현되고 각 지형에 따라 서로 다른 벡터들의 개수를 가진다.

제안한 알고리즘에서 조정해주어야 할 파라메터는 경쟁학습 신경회로망의 경우 식 (32)와 같이 놓았으며 Langevine 경쟁학습의 경우 식 (33)과 같이 놓았다. 초기 Weight vector 설정은 각 클래스 당 같은 개수의 데이터 벡터들로 이루어 지도록 했으며 Weight vector 수는 20개이다. 실험환경은 Pentium-100 IBM-PC에서 Visual C++ 1.5 컴파일러로 행하였으며 MFC함수를 사용, Windows Application으로 프로그램 하였다.

$$\epsilon(t) = 0.9 \cdot \left(1 - \frac{t}{\text{Number of Iteration}}\right) \quad (32)$$

$$\epsilon = 2^{-4} = \frac{1}{16} = 0.0625 \quad (33)$$

$$\lambda = 2^{-5} = \frac{1}{32} = 0.03125$$

식 (33)에서  $\epsilon$ 은 시간에 불변인 학습계수이며  $\lambda$ 는 이전 강화함수항의 아득이다.

표 1. FLC1 4-클래스 데이터에 대한 패턴인식 실험결과

Table 1. Experimental results of pattern recognition for FLC1 4 class Data.

| Epoch | 경쟁학습          |              | Langevine 경쟁학습 |              |
|-------|---------------|--------------|----------------|--------------|
|       | 학습 (%)        | 테스트 (%)      | 학습 (%)         | 테스트 (%)      |
| 100   | 93.525        | 92.09        | 94.625         | 92.93        |
| 200   | 94.0          | 90.6         | 96.125         | 94.40        |
| 300   | 94.25         | 93.0         | 94.75          | 93.73        |
| 400   | 94.625        | 91.13        | 95.875         | 93.53        |
| 500   | 94.625        | 92.4         | 96.125         | 94.87        |
| 600   | 95.875        | 91.93        | 95.875         | 94.73        |
| 700   | 94.625        | 88.47        | 96.375         | 94.60        |
| 800   | 95.375        | 88.0         | 96.125         | 94.20        |
| 900   | 94.75         | 83.13        | 95.0           | 93.93        |
| 1000  | 94.625        | 92.27        | 94.25          | 94.07        |
| 평균    | <b>94.625</b> | <b>92.07</b> | <b>95.51</b>   | <b>94.10</b> |

표 2. FLC1 8-클래스 데이터에 대한 패턴인식 실험결과

Table 2. Experimental results of pattern recognition for FLC1 8 class Data.

| Epoch | 경쟁학습         |              | Langevine 경쟁학습 |              |
|-------|--------------|--------------|----------------|--------------|
|       | 학습 (%)       | 테스트 (%)      | 학습 (%)         | 테스트 (%)      |
| 100   | 89.88        | 86.67        | 90.06          | 89.07        |
| 200   | 91.125       | 88.33        | 89.25          | 86.93        |
| 300   | 90.8125      | 88.0         | 90.06          | 88.4         |
| 400   | 89.625       | 86.57        | 90.0           | 88.1         |
| 500   | 90.6875      | 88.6         | 90.38          | 88.77        |
| 600   | 90.3125      | 86.80        | 90.50          | 89.07        |
| 700   | 89.625       | 86.57        | 90.88          | 89.17        |
| 800   | 89.625       | 86.57        | 90.50          | 89.43        |
| 900   | 89.9375      | 86.7         | 90.31          | 89.3         |
| 1000  | 89.88        | 86.67        | 90.13          | 88.5         |
| 평균    | <b>90.15</b> | <b>87.15</b> | <b>90.21</b>   | <b>88.57</b> |

표 3. Colorado 13-클래스 데이터에 대한 패턴인식 실험결과

Table 3. Experimental results of pattern recognition for Colorado 13 class Data.

| Epoch | 경쟁학습         |              | Langevine 경쟁학습 |              |
|-------|--------------|--------------|----------------|--------------|
|       | 학습 (%)       | 테스트 (%)      | 학습 (%)         | 테스트 (%)      |
| 100   | 48.31        | 49.46        | 62.40          | 61.42        |
| 200   | 47.52        | 48.66        | 62.10          | 61.23        |
| 300   | 48.31        | 49.46        | 62.0           | 61.33        |
| 400   | 48.31        | 49.36        | 60.22          | 60.04        |
| 500   | 48.31        | 49.46        | 63.49          | 62.61        |
| 600   | 47.52        | 48.67        | 62.80          | 61.82        |
| 700   | 48.31        | 49.46        | 61.51          | 60.53        |
| 800   | 48.31        | 49.46        | 59.52          | 59.94        |
| 900   | 47.52        | 48.67        | 60.81          | 59.45        |
| 1000  | 48.31        | 49.46        | 62.80          | 62.12        |
| 평균    | <b>48.31</b> | <b>49.21</b> | <b>61.77</b>   | <b>61.55</b> |

실험결과 제안한 알고리즘이 학습 및 테스트 양면에서 기존의 경쟁학습 신경회로망 알고리즘보다 평균적으로 더 나은 성능을 보였다. 신경회로망에 의한 추정 확률 분포와 실제 확률분포가 근사할 경우(FLC1 데이터), 즉 성능지표의 Convexity가 클 경우, 기존의 알고리즘보다 크게 분류 정확도가 향상되지 않으나, 성능지표의 Convexity가 멀어질 경우 (Colorado 13 클래스

데이터) Corollary 2의 결론을 통해 어느정도의 국소 최소점(Locally Optimal point)은 Langevine 경제학습 신경회로망이 극복 할 수 있음을 알 수 있다.

## VII. 결 론

본 논문에서는 먼저 Langevine 경제학습 신경회로망의 학습 알고리즘이 기존의 경제학습 알고리즘에 Gaussian 분포를 가진 Brownian Motion Process가 결합된 것임을 밝히고 이를 통해 Langevine 경제학습 신경회로망의 학습 알고리즘의 Fokker-Plank 방정식을 유도하였다. 유도된 Fokker-Plank 방정식은 Markov 준균위에 정의된 2차 편 미분 방정식이므로 선형 준균동형사상을 도입하여 Transition 확률의 축소(Contraction)가 어떤 조건 위에 성립하는지를 밝혀 대역 최소성의 조건을 밝혔으며, 실험을 통해 해석 결과를 살펴 보고 증명된 해석 조건들의 타당성을 실험을 통해 검증하였다.

본 논문에서 유도된 단위 단체(Simplex)들에 대한 대역 최소화 충분조건들은, 그러나 상당히 강한 조건들로서  $\int_{\Delta} L ds \rightarrow -\infty$ 가 만족 되려면  $t \rightarrow \infty$ 임을 의미 한다. 따라서 현재까지 알려진 거의 유일한 대역 최소화 알고리즘이 "Simulated Annealing"의 경우에서와 마찬가지로 유한번의 Epoch로 대역 최소점을 찾아내기는 사실상 어렵다는 것을 의미한다. 앞으로의 과제로는 본 논문에서 논한 분석이 매우 약한 위상에서의 분석임에 따라 좀 더 강한 위상에서, 전통적인 확률 미분 방정식의 해석방법으로 재 해석하여 본 알고리즘이 구체적으로 어떤 Nonconvex 조건까지를 만족하면서 대역 최소점에 도달 할 수 있는지를 밝히는 것과, Homology 분석을 통해 초기 Weight Vector들

의 입력공간 불활에 대한 해석을 하는 것이 요구된다.

## 참 고 문 현

- [1] J. Seok, S. Cho, "Self-Organizing Feature Map with Binary Reinforcement and Constant Adaptation Gain : For an easier Hardware Implementation", *Proc. ICONIP '94*, vol 2, 966-971 1, 1994.
- [2] I.I. Gihman, A.V. Skorohod, *The theory of stochastic process I*, Springer-Verlag, 1974.
- [3] H. Ritter, K. Schulten, "Convergence Properties of Kohonen's Topology Conserving Maps : Fluctuations, Stability, and Dimension Selection", *Biological Cybernetics*, vol 60, pp. 59-71, 1988.
- [4] I. I. Gihman, A. V. Skorohod, *The Theory of Stochastic Process II*, Springer-Verlag, 1974.
- [5] Yu. Ermoliev, R. J-B Wets (Eds.), 'Numerical Techniques for Stochastic Optimization', Springer-Verlag, 1980.
- [6] 석 진우, 조 성원, 홍 윤광, "시불변 학습계수와 이진 강화함수를 가진 경제학습 신경회로망에 의한 고장진단 폐지 전문가 시스템의 자동구축", *Proc '96 FAN 축제중합학술대회* 논문집, pp 153-158, 1996
- [7] T. S. Chiang, C. R. Hwang, S. J. Sheu, "Diffusion for Global Optimization in  $R^n$ ", *SIAM J. Cont. and Optimization*, vol 25, no 3, pp 737-753, May 1987,
- [8] P. Billingsley, *Probability and Measure*, Wiley, 3rd ED., 1994.

## 저 자 소 개

### 石 鎮 基(正會員)

1969년 6월 26일생. 1993년 2월 홍익대학교 전기제어 공학과 졸업.  
1995년 동대학원 전기 공학과 졸업  
(석사). 1997년 7월 현재 홍익대학교 대학원 전기 공학과 박사과정 재학중.  
주관심분야는 신경회로망, 확률과정,  
비선형 최적화, 비선형제어, 미분기하, 미분위상동임.

### 趙 成 元(正會員) 第32卷 B編 第1號 參照

