

Neurodynamika statystyczna głębokich sieci: geometria przestrzeni sygnałów

Mateusz Przyborowski

16 października 2020 r.

Główne wyniki z pracy:

1. Tensor metryczny, poprzez kolejne warstwy, zmienia się w sposób konforemny.
2. Krzywizna rozmaitości sygnału zbiega do stałej (pod pewnymi założeniami).
3. Wyjaśniony jest opis zmiany odległości pomiędzy dwoma sygnałami w kolejnych warstwach.

Oznaczenia

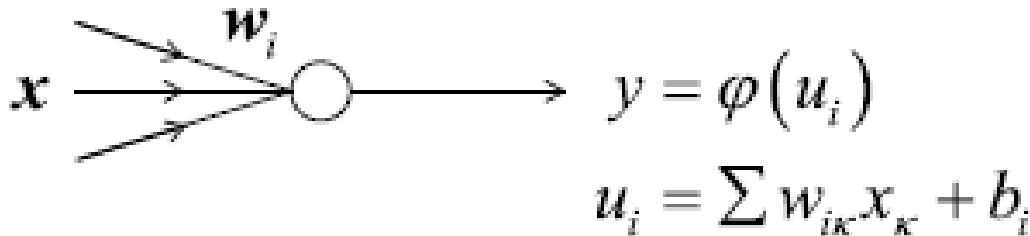
Rozważmy sieć neuronową o m neuronach i n -wymiarowym wejściu. Oznaczmy przez \mathbf{x} wejście, \mathbf{y} wynik sieci, w_{ik} waga połączenia pomiędzy k -tym wejściem i i -tym neuronem, b_i bias i -tego neuronu. Φ to funkcja błędu.

$$\mathbf{x} = (x_1, \dots, x_n)$$

$$\mathbf{y} = (y_1, \dots, y_m)$$

$$u_i = \sum_{\kappa} w_{i\kappa} x_{\kappa} + b_i$$

$$y_i = \varphi(u_i)$$



$$\varphi(u) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^u \exp\left\{-\frac{v^2}{2}\right\} dv$$

Oznaczenia

Zakładamy, że wagi połączeń są niezależnymi zmiennymi losowymi o rozkładzie normalnym ze średnią 0 oraz wariancjami σ^2/n i σ_b^2 dla wag i biasu. Wówczas u_i są niezależnymi zmiennymi losowymi wspólnego rozkładu normalnego o średniej 0 i wariancji: $\tau^2 = \frac{\sigma^2}{n} \sum x_k^2 + \sigma_b^2$

Dalej będziemy badali głęboką sieć neuronową, gdzie t -ta warstwa ma wejście x^{t-1} , które jest wynikiem działania $(t-1)$ -ej warstwy: $x^t = \varphi(W^t x^{t-1} + b^t)$

Podobnie zakładamy, że wagi i bias w_{ij}^t i b_i^t są niezależnymi zmiennymi losowymi rozkładów normalnych o średniej 0 i wariancjami odpowiednio σ_t^2/n_{t-1} i σ_b^2 .

X^t to przestrzeń możliwych wyników działania sieci w t -tej warstwie.

Przyjmując początkowy sygnał $X = X^0$, przez X^t oznaczmy jego obraz w t -tej warstwie.

Przetwarzanie sygnału wejściowego

Przy zadanym sygnale \mathbf{x}^t rozważmy funkcję ilości aktywności, zdefiniowaną jako:

$$A^t = \frac{1}{n^t} \sum (x_i^t)^2 = \frac{1}{n^t} \sum \{\varphi(u_i^t)\}^2 = E \left[\{\varphi(u_i^t)\}^2 \right]$$

By badać wielkości \mathbf{A}^t położmy również:

$$\chi_0(\sigma^2, A) = \int \varphi^2(\sigma_A v) Dv$$

$$\sigma_A^2 = \sigma^2 A + \sigma_b^2$$

$$Dv = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{v^2}{2}\right\} dv$$

$$\chi_1(\sigma^2, A) = \sigma^2 \int \{\varphi'(\sigma_A v)\}^2 Dv$$

Jako że \mathbf{u}_i^t ma rozkład normalny ze średnią 0 i wariancją $\tau_t^2 = \sigma_t^2 A^{t-1} + \sigma_b^2$:

$$\text{stąd mamy: } \mathbb{E} \left[\{\varphi(u_i^t)\}^2 \right] = \frac{1}{\sqrt{2\pi\tau_t}} \int \{\varphi(u)\}^2 \exp \left\{ -\frac{u^2}{2\tau_t^2} \right\} du = \chi_0(\sigma_t^2, A^{t-1})$$

W szczególności, przy funkcji błędu Φ zachodzi $\chi_0(\sigma^2, A) = \frac{1}{2\pi} \cos^{-1} \left(\frac{-\sigma_A^2}{1 + \sigma_A^2} \right)$

Twierdzenie.

Aktywność sygnału w sieci rozwija się zgodnie ze wzorem: $A^t = \chi_0(\sigma_t^2, A^{t-1})$

Jako że χ_0 jest monotonicznie rosnącą funkcją \mathbf{A} , gdy dla wszystkich \mathbf{t} mamy $\sigma_t^2 = \sigma^2$, to istnieje dokładnie jedno stabilne ekwilibrium $\hat{\mathbf{A}}$, a dla dużych \mathbf{t} \mathbf{A} zbiega do $\hat{\mathbf{A}}$.

Zdefiniowanie metryki

Niech $\mathbf{X} = \mathbf{X}^0$ będzie przestrzenią euklidesową. Połóżmy $dx = \sum dx_\kappa e_\kappa$ przy wektorach bazy ortonormalnej \mathbf{e} .

W kolejnych warstwach, proste z \mathbf{X} stają się krzywymi w $\hat{\mathbf{X}}^t$, stąd dla macierzy Jacobiego \mathbf{B} , $B_{ij}^t = \frac{\partial \bar{x}_i^t}{\partial \bar{x}_j^{t-1}} = \frac{\partial \varphi(u_i^t)}{\partial \bar{x}_j^{t-1}}$, mamy $d\bar{x}_i^t = \sum B_{ij}^t d\bar{x}_j^{t-1}$

Wektory bazowe przestrzeni stycznej do $\hat{\mathbf{X}}^t$ są zdefiniowane jako:

$$e_{\kappa i}^t = \frac{\partial \bar{x}_i^t}{\partial x_\kappa} = \sum_j B_{ij}^t e_{\kappa j}^{t-1}, \quad i = 1, \dots, n_t$$

Stąd wprowadzamy tensor metryczny: $g_{\kappa\lambda}^t = \langle \bar{e}_\kappa^t, \bar{e}_\lambda^t \rangle = \sum e_{\kappa i}^t e_{\lambda i}^t$

Zauważmy, że $\langle \tilde{e}_\kappa^t, \tilde{e}_\lambda^t \rangle = \sum_{i,j,k,l} B_{ik}^t B_{jl}^t \delta_{ij} e_{\kappa l}^{t-1} e_{\lambda k}^{t-1}$

Dla dość dużych n_t , z prawa wielkich liczb, mamy:

$$\frac{1}{n_t} \sum_i B_{ik}^t B_{il}^t = \mathbb{E} \left[\{\varphi'(u_i^t)\}^2 w_{ik}^t w_{il}^t \right] = \mathbb{E} \left[\{\varphi'(u_i^t)\}^2 \right] \mathbb{E} [w_{ik}^t w_{il}^t]$$

Natomiast z poprzednich wyników mamy:

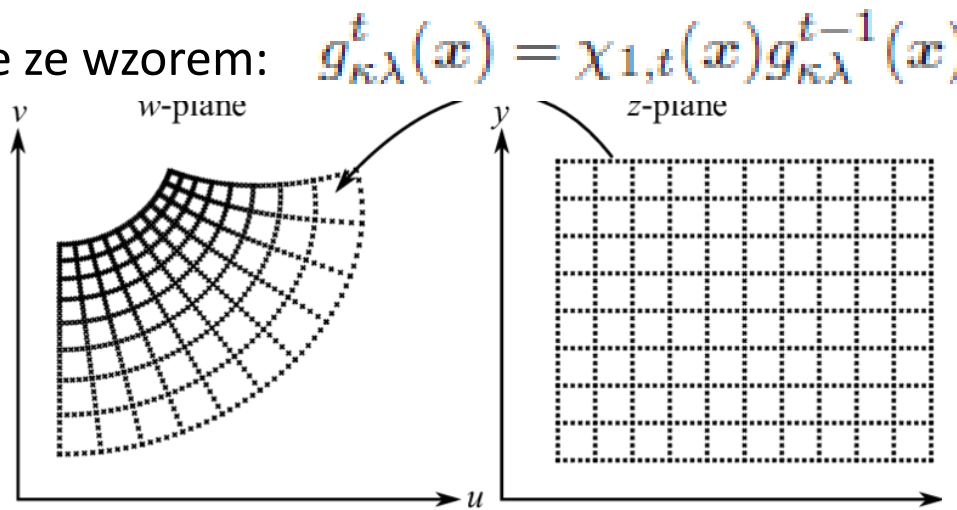
$$\chi_{1,t}(\sigma_t^2, A^{t-1}) = \sigma_t^2 \mathbb{E} \left[\{\varphi'(u_i^t)\}^2 \right]$$

$$\mathbb{E} [w_{ik}^t w_{il}^t] = \frac{\sigma_t^2}{n_t} \delta_{kl}$$

Twierdzenie.

Metryka jest przekształcana zgodnie ze wzorem: $g_{\kappa\lambda}^t(\mathbf{x}) = \chi_{1,t}(\mathbf{x}) g_{\kappa\lambda}^{t-1}(\mathbf{x})$

Jest to przekształcenie konforemne.



Metryka w warstwie \mathbf{t} -tej jest wyprowadzona jako $g_{\kappa\lambda}^t(\mathbf{x}) = \delta_{\kappa\lambda} \prod_{s=1}^{t-1} \chi_{1,s}$

Dla $\sigma_t^2 = \sigma^2$ i dla dość dużych \mathbf{t} , \mathbf{A}^t zbiega do \mathbf{A} . Stąd asymptotycznie mamy:

$$g_{\kappa\lambda}^t(\mathbf{x}) \approx \{\bar{\chi}_1\}^{t-1} \delta_{\kappa\lambda}, \quad \bar{\chi}_1 = \chi_1(\sigma^2, \bar{A})$$

Dla ϕ funkcji błędu mamy $\chi_1(\sigma^2, A) = \frac{\sigma^2}{2\pi} \frac{\sigma_b^2 + \sigma^2 A}{\sqrt{1 + 2(\sigma_b^2 + \sigma^2 A)}}$

Gdy $\bar{\chi}_1 < 1$, $g_{\kappa\lambda}^t$ zbiega do zera.

Gdy $\bar{\chi}_1 > 1$ (dzieje się tak przy dużej wariancji), długość krzywych rośnie wraz ze wzrostem \mathbf{t} , ale skoro sygnał leży w ograniczonej podprzestrzeni \mathbf{X}^t (z wyjątkiem niektórych funkcji aktywacji, np. ReLU), to krzywa musi być bardzo poskręcana.

Krzywizna

Krzywizna rozmaitości \mathbf{X}^t jest mierzona przez to, jak bardzo bazowe wektory $\hat{\mathbf{e}}_k^t$ przestrzeni stycznej do \mathbf{X}^t się zmieniają gdy punkty poruszają się w kierunku $\hat{\mathbf{e}}_l^t$. Za pomocą pochodnej kierunkowej kładziemy $H_{\kappa\lambda}^t = \nabla_\lambda \tilde{\mathbf{e}}_\kappa^t$

Zauważmy, że dla $\partial_\kappa = \frac{\partial}{\partial x_\kappa}$ mamy:

$$\begin{aligned} H_{\kappa\lambda i}^t &= \partial_\lambda \partial_\kappa \varphi(u_i^t) = \partial_\lambda e_{\kappa i}^t = \partial_\kappa \left\{ \sum_j \varphi'(u_i) w_{ij} \partial_\lambda \tilde{x}_j^{t-1} \right\} = \\ &= \varphi''(u_i) (w_i \cdot \partial_\kappa \tilde{x}^{t-1}) (w_i \cdot \partial_\lambda \tilde{x}^{t-1}) + \varphi'(u_i) w_i \cdot \partial_\kappa \partial_\lambda \tilde{x}^{t-1} = \\ &= \varphi''(u_i) (w_i \cdot \tilde{\mathbf{e}}_\kappa^{t-1}) (w_i \cdot \tilde{\mathbf{e}}_\lambda^{t-1}) + \varphi'(u_i) w_i \cdot \mathbf{H}_{\kappa\lambda}^{t-1} \end{aligned}$$

Przy ostatniej równości korzystamy z tego, że $\mathbf{H}_{\kappa\lambda}^{t-1} = \partial_\kappa \partial_\lambda \tilde{x}^{t-1}$

Teraz możemy położyć jako wielkość krzywizny $|\mathbf{H}_{\kappa\lambda}^t|^2 = \langle \mathbf{H}_{\kappa\lambda}^t, \mathbf{H}_{\kappa\lambda}^t \rangle = \sum_i (H_{\kappa\lambda i}^t)^2$

Z prawa wielkich liczb mamy:

$$\begin{aligned} |H_{\kappa\lambda}^t|^2 &= n_t E \left[\{\varphi''(u_i)\}^2 (w \cdot \tilde{e}_\kappa^{t-1})^2 (w \cdot \tilde{e}_\lambda^{t-1})^2 \right] \\ &\quad + 2n_t E \left[\varphi'(u_i) \varphi''(u_i) (w \cdot \tilde{e}_\kappa^{t-1}) (w \cdot \tilde{e}_\lambda^{t-1}) (w \cdot \partial_\kappa \tilde{e}_\lambda^{t-1}) \right] \\ &\quad + n_t E \left[\{\varphi'(u_i)\}^2 (w \cdot H_{\kappa\lambda}^{t-1})^2 \right] \end{aligned}$$

Środkowy wyraz upraszcza się z nieparzystości funkcji. Pozostałe dwa można przeformułować do postaci:

$$|H_{\kappa\lambda}^t|^2 = \chi_1(\sigma_t^2, A^{t-1}) |H_{\kappa\lambda}^{t-1}|^2 + \frac{1}{n_t} (1 + 2\delta_{\kappa\lambda}) \{\chi_2(\sigma_t^2, A^{t-1})\} (g_{\kappa\kappa}^{t-1})^2 (g_{\lambda\lambda}^{t-1})^2$$

Dalej upraszczając kładziemy:

$$\gamma_t^2 = \frac{1}{n_t} \sum (g_t^{-1})^{\kappa\mu} (g_t^{-1})^{\lambda\nu} H_{\kappa\lambda i}^t H_{\mu\nu j}^t \delta_{ij} = \frac{\gamma_{t-1}^2}{\chi_{1,t-1}} + 3 \frac{\chi_{2,t-1}}{n_t (\chi_{1,t-1})^2} = 3 \frac{1}{n_t} \sum_{s=1}^{t-1} \frac{\chi_{2,t-s}}{\chi_{1,t-s}^2} \left(\prod_{r=1}^s \frac{1}{\chi_{1,t-r}} \right)$$

Dla dość dużych t , $\sigma_t^2 = \sigma^2$ i $\bar{\chi} > 1$ mamy $\gamma_t^2 = \frac{3\chi_2}{n_t \chi_1 (\chi_1 - 1)}$,
co oznacza, że skalar krzywizny zbiega do stałej.

Dla $\bar{\chi}_1 < 1$ jest rozbieżny.

Odległość między sygnałami

Niech \mathbf{x}^{t-1} i \mathbf{y}^{t-1} (\mathbf{x}^t i \mathbf{y}^t) oznaczają dwa sygnały na wejściu (wyjściu) t -tej warstwy.

Położmy $D_t = D(\mathbf{x}^t, \mathbf{y}^t) = \frac{1}{n_t} \sum (x_i^t - y_i^t)^2$ i $C_t = C(\mathbf{x}, \mathbf{y}) = \frac{1}{n_t \sqrt{A(\mathbf{x}^t) A(\mathbf{y}^t)}} \sum x_i^t y_i^t$

do badania odpowiednio odległości i nachodzenia się sygnałów. Zachodzi:

$$D(\mathbf{x}, \mathbf{y}) = A(\mathbf{x}) + A(\mathbf{y}) - 2\sqrt{A(\mathbf{x})A(\mathbf{y})}C(\mathbf{x}, \mathbf{y})$$

A więc znając $C_t = \psi(C_{t-1})$ możemy wyrazić $D_t = \xi(D_{t-1})$

Założmy dodatkowo, że $A(\mathbf{x}) = A(\mathbf{y}) = A^t$; wówczas mamy $D_t = 2A^t(1 - C_t)$

a więc możemy wyrazić odległość jako $\xi(D_{t-1}) = 2A^t \left\{ 1 - \psi \left(1 - \frac{D_{t-1}}{2A^t} \right) \right\}$

Zmienne losowe $u = \sum w_j x_j + b_j$ i $u' = \sum w_j y_j + b_j$ mają rozkład łączny normalny o zerowej średniej i macierzy kowariancji:

$$E[u^2] = E[u'^2] = \sigma_{A^t}^2,$$

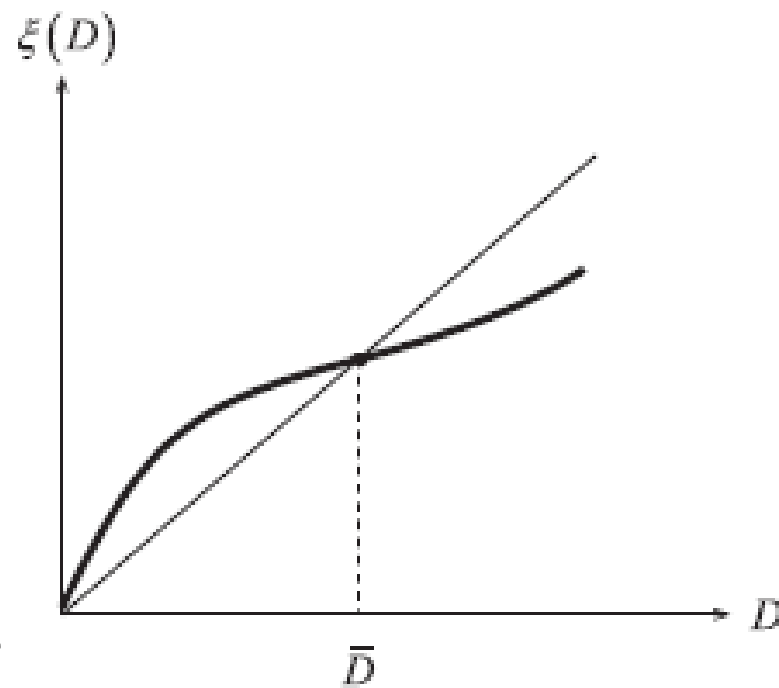
$$E[uu'] = \sigma^2 A^t C(x, y) + \sigma_b^2 = \sigma_{A^t C}^2.$$

Wówczas sygnały natchodzą się zgodnie z: $C = \frac{1}{A^t} E[\varphi(u)\varphi(u')]$
co daje nam ogólną postać:

$$C_t = \frac{1}{2\pi A^t} \cos^{-1} \left(-\frac{C_{t-1} A^t \sigma_t^2 + \sigma_b^2}{\sigma_{A^t}^2 \sqrt{1 + \sigma_{A^t}^2}} \right)$$

A z tego wyprowadzamy ξ . W szczególności ξ jest monotoniczną funkcją D i mamy $\xi(0) = 0$ oraz, przy $\sigma_t^2 = \sigma^2$, zachodzi $\xi'(0) = \bar{\chi}_1$.

Dla $\bar{\chi} > 1$ mamy jeszcze jedno, unikalne rozwiązanie $\bar{D} = \xi(\bar{D})$, które jest jedynym stabilnym przy t zbiegającym do nieskończoności. W tym miejscu teoria bazuje na nieprzeliczalnej kardynalności przestrzeni X .



Wnioski

1. Przy $\bar{\chi}$ bliskiej **1** dynamika \bar{x}^t jest bliska chaotycznej.
2. Przy $\bar{\chi}$ równej **1** krzywizna jest rozbieżna do nieskończoności (ale pod warunkiem że **n** skończone!)
3. Otrzymane wyniki opisują zachowanie przy $t \rightarrow \infty$ oraz $n \rightarrow \infty$, ale w praktyce wartości **t** oraz **n** pozostają skończone.
4. Autorzy podkreślają potrzebę badań wpływu „skończoności” na dotychczas uzyskane wyniki.

Dziękuję za uwagę

Bibliografia:

Shun-ichi Amari, Ryo Karakida, and Masafumi Oizumi.

Statistical Neurodynamics of Deep Networks: Geometry of Signal Spaces.

Technical report, 2018a.