

# GEOEMTRIA STRUCTURILOR REȚELELOR NEURALE

Vlad Raul Constantinescu<sup>1</sup>, Ionel Popescu<sup>2,3</sup>

<sup>1</sup> Școala doctorală a Facultății de Matematică, Universitatea din București  
vlad\_constantinescu95@yahoo.com

<sup>2</sup> Facultatea de Matematică și Informatică, Universitatea din București

<sup>3</sup> Institutul de Matematică al Academiei Române "Simion Stoilow"  
ioionel@gmail.com

O formulare matematică a problemei de învățare automată este următoarea. Presupunem că avem un set de date  $(X_i, Y_i)_{i=1, \dots, N}$  extrase independent după o distribuție  $\mu$ , unde  $X_i$  sunt datele de input iar  $Y_i$  sunt etichetele asociate. De exemplu, dacă avem un set de imagini, avem  $X_i$  este imaginea în sine iar  $Y_i$  este descrierea imaginii (câine, pisică, avion, etc). În acest context problema este aceea de a determina o funcție  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  în așa fel încât să avem  $Y_i = f(X_i)$ . Este un fel de problema de regresie. În mod realist este dificil de găsit o astfel de funcție datorită multor inconveniente, de exemplu, avem zgomot, sau o astfel de funcție este mult prea complicată. Așa că în general se caută o astfel de funcție într-o familie destul de simplă. Pe de altă parte, cum ziceam, nu ne putem aștepta ca  $f$  să descrie foarte bine datele, așa că trebuie să definim un fel de eroare asociată. În acest context se definește o funcție de pierdere,  $L(Y', Y)$  care măsoară cât de departe este  $Y'$  de  $Y$ , unde  $Y = f(X)$  este predicția și  $Y'$  este eticheta reală a lui  $X$ . O astfel de funcție la îndemână este  $L(Y', Y) = \|Y - Y'\|^2$ , însă sunt și alte funcții care joacă un rol mai bun în astfel de situații, cum ar fi entropia relativă pentru cazul în care  $Y, Y'$  iau doar valori discrete și sunt probabilități.

Rețelele neurale standard determină o parametrizare a acestor funcții dată de structura rețelei. Putem interpreta aceasta ca o funcție de la un spațiu  $\mathbb{R}^k$  în spațiul de funcții  $\mathcal{F}(\mathbb{R}^n; \mathbb{R}^m)$ . Notăm această funcție cu  $V : \mathbb{R}^k \rightarrow \mathcal{F}(\mathbb{R}^n; \mathbb{R}^m)$ . Pentru o pereche  $(X, Y) \in \mathbb{R}^n \times \mathbb{R}^m$  și folosind funcția  $L$  de mai sus, putem defini acum  $F_{(X, Y)}(\theta) = L(V(\theta)(X), Y)$  care ia valori reale. Având la îndemână un set de date  $(X_i, Y_i)_{i=1, \dots, N}$ , în principiu vrem să optimizăm (minimizăm mai precis) după funcția pierdere medie de forma

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N F_{(X_i, Y_i)}(\theta).$$

În principiu noi am dori de fapt minimizarea riscului mediu data de  $\mathbb{E}[F_{(X, Y)}(\theta)]$ , unde media este calculată în raport cu măsura de probabilitate  $\mu$  ce generează datele. Din păcate această distribuție nu este cunoscută, însă putem aproxima această medie cu varianta discretă de mai sus.

Metoda standard de a găsi  $\theta$  optim pentru funcția  $\mathcal{L}$  este bazată pe metoda *gradient descent*. Această metodă constă în găsirea unei linii poligonale cu direcția opusă gradientului lui  $\mathcal{L}$  care să ne conducă spre un minim local. Asta implică așadar calculul gradientului lui  $\mathcal{L}$  și în particular pentru fiecare  $F_{(X_i, Y_i)}$ . Acesta din păcate este foarte costisitor, așa că însăși problema calculului gradientului lui  $\mathcal{L}$  devine intractabil. Modificarea standard aplicată este cea de mini batch, adică în loc să calculăm la fiecare pas  $\mathcal{L}$  și gradientul asociat, calculăm gradientul doar de-a lungul unui subsample mai mic din datele disponibile. Așa că avem acum de calculat doar un număr mai mic de date. Aceasta metoda este cunoscută ca metoda *stochastic gradient descent* pentru că alegem de fiecare dată un subsample mai mic și la fiecare pas ne mișcăm în direcția acestui gradient generat de subsample.

Pentru rețelele neurale, încă nu există o justificare suficientă pentru marele mister, anume faptul că aceste metode sunt foarte eficiente din punct de vedere practic.

Una din problemele pe care le propunem pentru această teză este studierea proprietăților acestor funcții de risc  $F_{(X, Y)}(\theta)$ . În particular am dori să vedem ce proprietăți intrinsece trebuie să aibă aceste clase de funcții ce garantează metoda gradientului să funcționeze așa de bine.

Una din direcțiile de studiu este legată de flexibilizarea structurilor de rețele neurale și studiul mai întâi empiric ale acestor rețele de tip neural. De exemplu, putem în loc de rețele de tip standard să gândim niște rețele cu o topologie de tip diferit. Mai exact spus, rețelele clasice sunt de tipul stratul 0 citește datele și următorul strat este compus din termeni de stratul strict de dinainte. În acest fel topologia este oarecum liniară. Există și alte structuri ceva mai complexe, în care de exemplu straturile interacționează între ele într-un anumit fel. Ce ne propunem aici este analiza structurilor și topologiei asociate. Un exemplu simplu în acest caz ar fi cel în care straturile sunt dispuse pe un tor și rețeaua neurală este organizată pe un mesh pe tor. În acest fel stratul de input și de output constând într-un număr de noduri alese aleator pe acest mesh.

O altă direcție de lucru este cea în care vrem să investigăm ce proprietăți, sa le numim  $P$ , sunt importante pentru aceste rețele neurale. De exemplu, o primă abordare ar fi cea în care să studiem ce proprietăți se păstrează în restricțiile acestor funcții de tip  $F_{(X, Y)}$  de-a lungul unor drepte și plane luate aleator în  $\mathbb{R}^k$  (spațiul ambiental al parametrilor). Este foarte posibil ca acest tip de structură aleatoare văzută doar la nivel de restricții pe structuri geometrice simple să aibă ceva în comun. Următorul pas în această analiză ar fi cel legat de proprietățile ce se păstrează la adunare, adică date fiind două astfel de funcții cu proprietatea  $P$ , suma lor să păstreze aceste proprietăți. O astfel de proprietate este foarte importantă deoarece ne poate furniza diverse alte structuri potențial generale care să reproducă rețelele neurale înșa într-o formă mult mai simplă și în același timp cu o justificare matematică riguroasă.

## Bibliografie

- [1] Omer Bobrowski, Matthew Kahle, Primoz Skraba, et al. Maximally persistent cycles in random geometric complexes. *The Annals of Applied Probability*, 27(4):2032–2060, 2017
- [2] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT press Cambridge, 2016

- [3] Qianli Liao and Tomaso Poggio. Theory ii: Landscape of the empirical risk in deep learning. *arXiv preprint arXiv:1703.09833*, 2017.
- [4] Tomaso Poggio, Kenji Kawaguchi, Qianli Liao, Brando Miranda, Lorenzo Rosasco, Xavier Boix, Jack Hidary, and Hrushikesh Mhaskar. Theory of deep learning iii: explaining the non-overfitting puzzle. *arXiv preprint arXiv:1801.00173*, 2017.
- [5] Christian Robert. Machine learning, a probabilistic perspective, 2014.
- [6] Jurgen Schmidhuber. Deep learning in neural networks: An overview. *Neural networks*, 61:85–117, 2015.
- [7] Herut Uzan, Shira Sardi, Amir Goldental, Roni Vardi, and Ido Kanter. Biological learning curves outperform existing ones in artificial intelligence algorithms. *Scientific reports*, 9(1):1–11, 2019.
- [8] Saining Xie, Alexander Kirillov, Ross Girshick, and Kaiming He. Exploring randomly wired neural networks for image recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1284–1293, 2019.
- [9] Chiyuan Zhang, Qianli Liao, Alexander Rakhlin, Brando Miranda, Noah Golowich, and Tomaso Poggio. Theory of deep learning iib: Optimization properties of sgd. *arXiv preprint arXiv:1801.02254*, 2018.