

UWAGI od Rafala  
(rafjoz@gmail.com)

Wydaje mi się że może być warto poeksperymentować z inną parametryzacją predykcji  $p(y | x)$ . W sieciach neuronowych optymalizujemy zwykle  $-\log p(y | x)$ , gdzie  $p$  jest modelowane jakimś rozkładem prawdopodobieństwa. Jest tutaj spore pole do wyboru, w zależności od rodzaju danych.

Np. gdy  $p$  jest modelowane rozkładem Gaussowskim ze stałą wariancją ( $p \sim N(y | x; \text{std}=1)$ ), odpowiadająca funkcja kosztu będzie:

$-\log N(y | x; \text{std}=1) = -\log(1/\sqrt{2 * \pi * \text{std}^2} * e^{-(y-y')^2 / 2\text{std}^2}) = 0.5 * (y-y')^2 + \text{const} == \text{MSE}(y, y')$ , gdzie  $y'$  to predykcja sieci neuronowej

Z powyższego wynika że używając funkcji kosztu MSE, efektywnie zakładamy że nasze etykiety są dobrze modelowane rozkładem normalnym. Te założenie może nie być idealne, jako że wiemy, iż prawdopodobieństwa katów mogą mieć 2 wartości maksymalne w niektórych przypadkach – wtedy sieć neuronowa zwraca najgorszy możliwy wynik, tzn. wartość po środku między dwoma maksimumami.

Można ten problem rozwiązać na kilka sposobów modelując predykcje innymi rozkładami prawdopodobieństwa.

1. Najprostszy z nich, który był testowany w tej pracy to podzielenie dziedziny funkcji  $[0;1]$  na wiele kubelków i modelowanie wyników rozkładem dyskretnym. Rozkład ten jest dużo bardziej elastyczny i pozwala modelować o wiele bardziej skomplikowane rozkłady, ale z drugiej strony przyjmuje dużo mniej założeń o strukturze modelowanego rozkładu. W naszym przypadku wiemy, że nasza funkcja jest ciągła, więc predykcje dla sąsiednich wartości powinny być podobne. Prawdopodobnie z tego powodu może być trudniejsze trenowanie takiej sieci neuronowej.

(\*) Można spróbować z dodaniem jakichś regularyzacji do sieci, która pomoże w treningu, np. dodatkowa funkcja kosztu, która mówi że sąsiednie predykcje muszą być podobne, tzn.:

$\text{jakis\_wspolczynnik} * \sum (p_i - p_{i+1})^2$

2. Inny sposób to użycie mixture of gaussians jako funkcji kosztu, tzn.

$-\log[p_1 * N_1(y_1 | x; \text{std}=1) + p_2 * N_2(y_2 | x; \text{std}=1)]$ ,  $p_1 + p_2 = 1$

Taka funkcja kosztu pozwoli nam modelować nasz rozkład w lepszy sposób, bo będziemy w stanie przewidzieć potencjalne oba maksima rozkładu prawdopodobieństwa. W tym wypadku sieć neuronowa będzie zwracała 4 wartości zamiast 1, tzn.  $[p_1, p_2, y_1, y_2]$  do parametryzacji obu funkcji Gaussowskich, a nasza funkcja kosztu będzie trochę bardziej skomplikowana

(ale wciąż dość prosta do optymalizacji używając numerycznie stabilnej operacji  $\text{logsumexp}()$ ):

$-\log p_1 N_1 + p_2 N_2 = -\log[p_1 / \sqrt{2\pi} * e^{-(y-y_1)^2/2} + p_2 / \sqrt{2\pi} * e^{-(y-y_2)^2/2}] = -\log[\sum_i e^{\log(p_i) + \log(1/\sqrt{2\pi}) - (y-y_i)^2/2}] = -\text{logsumexp}[\log(p) - (y-y')^2/2]$ , gdzie  $p=[p_1, p_2]$ ,  $y'=[y_1, y_2]$

Tutaj, jest przykładowa implementacja w tensorflow: <https://ngbinghao.gitlab.io/posts/mixture-density-networks-basics/> dla trochę ogólniejszego przypadku, w którym dodatkowo jeszcze parametryzujemy odchylenia standardowe gaussów sieciami neuronowymi, tzn. sieć zwraca  $p=[p_1, p_2]$ ,  $y'=[y_1, y_2]$  i  $\text{std}=[\sigma_1, \sigma_2]$ , co odpowiada rozkładowi:  $p_1 * N(\text{mean}=y_1, \text{std}=\sigma_1 | x) + p_2 * N(\text{mean}=y_2, \text{std}=\sigma_2 | x)$

Funkcja kosztu jest wtedy podobna, ale dochodzi jeszcze kilka dodatkowych składników (od  $\sqrt{2\pi} \sigma^2$  które już nie jest stałą wartością) i  $((y-y')/\sigma)^2/2$

Ta implementacja jest całkiem dobra, ale niestety ma jeden błąd numeryczny,  $\log(p)$  powinno być wyliczone używając funkcji `tf.nn.log_softmax(p)` zamiast `tf.log(tf.nn.softmax(p))`. Ta druga może nie być stabilna numerycznie.

W tym wariancie, podobnie jak w użytym w raporcie, optymalizację można wykonać względem wartości maksymalnej. Można też spróbować czegoś bardziej skomplikowanego, co wzięłoby pod uwagę cały rozkład zamiast tylko maksimum, ale nie wiem czy warto próbować. Byłoby to trochę trudniejsze chyba i może nie pomagać.

(\*) Znalazienie maksimum wymaga tutaj sprawdzenia wartości dla 0, 1 i w punktach  $y_1$  i  $y_2$  (przy założeniu że są w przedziale  $[0;1]$ ) jako że maksymalna wartość może być w jednym z tych 4 punktów w zależności od predykcji sieci.

3. Warianty 1. i 2. z parametryzowanym odchyleniem standardowym rozkładów też mogłoby być ciekawe do sprawdzenia, ale prawdopodobnie jest mniej ważne w tym problemie

4. Jeden problem o którym nie wspomniałem wcześniej jest taki że warianty z rozkładami normalnymi nie biorą pod uwagę tego że dziedziną  $y$  w naszym problemie jest zamknięty przedział  $[0;1]$  podczas gdy dziedziną rozkładów normalnych jest całe  $\mathbb{R}$ . Można to spróbować użyć i było kilka ciekawych prac na ten temat w przypadku modeli które generują obrazy. Tzn intensywności kolorów pikseli są modelowane 3 zmiennymi o wartościach z przedziału  $[0;255]$  dla każdego z kanałów RGB.

Historycznie implementacje które używały rozkładu normalnego do modelowania tych wartości dawało lepsze rezultaty niż rozkłady dyskretnie. Jakis czas temu ludzie zaczęli eksplorować inne rozkłady które starają się wziąć pod uwagę informacje na temat tego że dziedziną funkcji jest ograniczona. Jedną z prac jest: <https://arxiv.org/abs/1701.05517> która używa mikstury zdyskretyzowanych rozkładów logistycznych, co wydaje się działać bardzo dobrze na obrazkach.

Jest nawet kod załączony do tej pracy tutaj: <https://github.com/openai/pixel-cnn>. Nie jestem pewien czy warto to próbować jako że może być to dość skomplikowane, ale wciąż wydaje się być ciekawe i może dać lepsze wyniki.

Inne uwagi:

\* Sposób ewaluacji wydaje mi się że mógłby być lepszy, tzn odległość od maksimum nie jest idealna w sytuacji gdy rzeczywisty rozkład ma dwa maksima i w takich wypadkach odległość może być niejednoznaczna gdy mamy 2 punkty o takich samych wartościach. Dodatkowo, wydaje mi się że ważniejsze/ciekawsze jest modelowanie całego rozkładu niż tylko wartości maksymalnej, szczególnie w sytuacjach gdzie skalar/pseudoskalar (0 i 1) mają podobne wartości a maksimum lub minimum jest gdzieś po środku. Tutaj ewaluacja by użyła punktów 0, 0.5 lub 1 jako optymalnej i ucieka tu bardzo dużo interesujących informacji o rzeczywistym rozkładzie.

\* Inny sposób mógłby po prostu ewaluować całą entropię krzyżową między rzeczywistym rozkładem a przewidzianym. W przypadku rozkładu dyskretnego efektywnie to robimy, ale takiej samej metryki można użyć we wszystkich innych rozkładach. Konkretnie, dla rozkładu normalnego, możemy zewaluować funkcję w 50 punktach na przedziale  $[0; 1]$ , wyliczyć prawdopodobieństwa, które trzeba później znormalizować żeby się sumowały do 1 (co efektywnie ogranicza dziedzinę predykcji do przedziału  $[0; 1]$ ) Można to pewnie zrobić porządniej używając całek jako że znamy funkcyjną reprezentację tych rozkładów ale może to nie być konieczne.

\* Komentarz do sekcji 6 ("Regresja do parametrów rozkładu wag A, B, C").

Wydaje mi się że mogłoby być użyteczne spróbować połączenie tej metody razem z poprzednimi, tzn używać tej regresji do A, B, C jako dodatkowej funkcji kosztu, przynajmniej na początku treningu. Może to przyspieszyć uczenie i potencjalnie dostarczyć sieci neuronowej dodatkowych informacji.

\* Jako że wiemy, że rzeczywisty rozkład używa funkcji  $\sin(x)$  i  $\cos(x)$ , może być warto użyć tych przekształceń jako nieliniowości w samej sieci neuronowej. Konkretnie, w ostatniej warstwie, którą przekształcamy nieliniowością ReLU, można dodatkowo przekształcić wartości  $\sin(x)$  i  $\cos(x)$  i skleić je ze sobą, co nam da 3x więcej wyjść z tej konkretnej warstwy:

```
x = wyjście poprzedniej warstwy
nowa_warstwa = tf.concat([tf.nn.relu(x), tf.sin(x), tf.cos(x)], axis=1) # <- nowe transformacje
y = linear(nowa_warstwa, 1)
loss = tf.losses.mean_squared_error(...)
```

Cieško powiedzieć czy to pomoże, ale jest bardzo proste do sprawdzenia i może być użyteczne.

Ach, jeszcze jedno - od strony praktycznej polecam przetestować numerycznie te nowe rozkłady prawdopodobieństwa, jako że łatwo popełnić błędy przy definiowaniu funkcji kosztu, ale bardzo łatwo sprawdzić, czy wartości tych funkcji są poprawne, używając najprostszych wzorów, które nie muszą być bardzo stabilne numerycznie. Np., dla rozkładu normalnego lub mikstury dwóch rozkładów, bardzo łatwo ewaluować  $N(y'; \text{mean}=x', \text{std}=1)$  używając Scipy lub innych bibliotek do obliczeń numerycznych.