

HW (6) باسمہ تعالیٰ

(12)

$$\text{Swish} = n \cdot \text{Sigmoid}(n) \quad (1)$$

میش = سواش پر سواش

$$\text{Mish} = n \tanh(\text{Softplus}(n))$$

میش = میس پر میس

(2)

$$\text{Swish} = n \cdot \text{Sigmoid}(n)$$

$$S'(n) = \text{Sigmoid}(n) + n \cdot \text{Sigmoid}(n)(1 - \text{Sigmoid}(n))$$

$$= \sigma(n) + n \cdot \sigma(n) - n \cdot \sigma(n)^2$$

$$= n \cdot \sigma(n) + \sigma(n)(1 - n \cdot \sigma(n))$$

$$= S(n) + \text{Sigmoid}(n)(1 - S(n))$$

سوال 1st-derivative-Swish problem

$$n(n) = n \cdot \tanh(\text{Softplus}(n))$$

$$\tanh(n) = \frac{e^n - e^{-n}}{e^n + e^{-n}}$$

$$\text{Softplus} = \ln(1 + e^x)$$

$$n(n) = x \cdot e^{\ln(1+e^n)} - e^{-\ln(1+e^n)} \Big/ \left(e^{\ln(1+e^n)} + e^{-\ln(1+e^n)} \right)$$

Der

$$n(n) = \frac{e^n - e^{-n}}{e^n + e^{-n}}$$

$$n'(n) = n \cdot \left(\frac{(e^n + 1)^2 - 1}{(e^n + 1)^2 + 1} \right)$$

سوال 2nd-derivative-mish problem

تابع Relu نسبت به Sigmoid و tanh

هزینه‌های حسابی کمتری دارد و در نتیجه سریع‌تر است

در تابع Sigmoid ، tanh به این دلیل که در

$+\infty$ و $-\infty$ نسبت به صفر را دارند در نتیجه همه مدل‌ها در نتیجه

یادگیری کند و ~~در نتیجه~~ ~~در نتیجه~~ ~~در نتیجه~~

تابع Relu در مقدارهای منفی مقدار صفر حلقه را

برمی‌گیرد و آنکه در نتیجه نسبت به صفر دارد و عمل ~~در~~ مقدار

منفی ای بدست آید مدل ~~در~~ در یادگیری دچار مشکل می‌شود

هر چند که روی یک batch لحاظ می‌شود و احتمال اینکه

همه مقدارهای منفی شوند بسیار کم است

دو عمل Mish ، Swish در مقادیر منفی

(۲۰)

ی

نسبت (۱) مقدار های گفته شده را نسبت به مقدار ها

Initial رتبه ها نسبت اگر مقدار ها کمتر Initial

شوند داریم \rightarrow epoch اول :

$$MSE = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

$$y \begin{cases} 0 & (0 - 0.5)^2 = 0.25 \\ 1 & (1 - 0.5)^2 = 0.25 \end{cases}$$

binary cross entropy: $-(y \log(\hat{y}) + (1-y) \log(1-\hat{y}))$

$$y \begin{cases} 0 & -(0 + \log(0.5)) \approx 0.7 \\ 1 & -(0 + \log(0.5)) \approx 0.7 \end{cases}$$

$$1 - (0 + \log(0.5)) \approx 0.7$$

ج ۱- چون مدل دیتاهای ℓ_2 حفظ کرده و overfit شده است

ج ۲- برای min cross entropy در epoch 60

قبل از آگرای

برای MSE در epoch 80 یا 100

ج ۳-

باتوجه به مقدار مقدار 0.5 از همه بهتر بود

چون مقدار Validation ^{loss} در 0.5 بسیار خوب به هم رسید

در نظری رسید overfit رخ نداده

در مقدار $\alpha = 1$ به جز در epoch آخر، اگر اند