

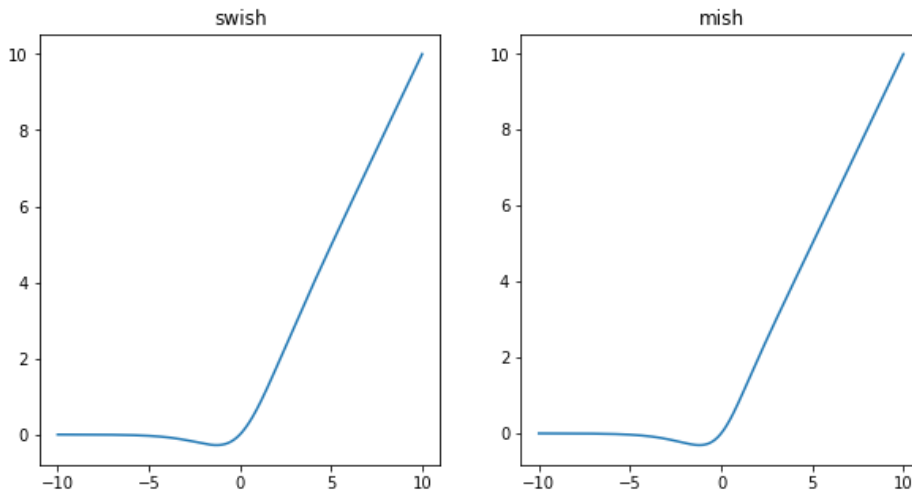
۱- الف) رابطه *swish* :

$$f(x) = x * \text{sigmoid}(x) = \frac{x}{1 + e^{-\beta x}}$$

رابطه *mish* :

$$f(x) = x * \tanh(\text{softplus}(x)) = x * \tanh(\ln(1 + e^x))$$

نمودارها:



ب) مشتق *swish* :

$$\begin{aligned} f'(x) &= x' \text{sig}(x) + x \text{sig}'(x) = \text{sig}(x) + x \text{sig}(x)(1 - \text{sig}(x)) \\ &= \text{sig}(x) + x \text{sig}(x) - x \text{sig}(x)^2 \\ &= x \text{sig}(x) + \text{sig}(x)(1 - x \text{sig}(x)) \\ &= f(x) + \text{sigmoid}(x)(1 - f(x)) \end{aligned}$$

مشتق *mish* :

$$\begin{aligned} f'(x) &= \tanh(\ln(1 + e^x)) + \frac{x e^x}{1 + e^x} \text{sech}^2(\ln(1 + e^x)) \\ &= \text{sech}^2(\text{softplus}(x)) * x * \text{sigmoid}(x) + \frac{f(x)}{x} \\ &= \text{sech}^2(\text{softplus}(x)) * \text{swish}(x) + \frac{f(x)}{x} \end{aligned}$$

ت) *relu* نسبت به *sigmoid* و *tanh* سریع تر است زیرا با مثبت بودن ورودی گرادین می‌تواند بیشتر شود و بیشتر تغییر کند. از طرف دیگر حد بالا ندارد و این باعث می‌شود که گرادین اشباع نشود.

در رابطه با *swish* و *mish* هر دو مانند *relu* از بالا نامحدود هستند و از پایین حد دارند، ولی به دلیل *non-monotonicity* هر دو مقادیر منفی کوچکی تولید میکنند که باعث جریان گرادین در مقادیر زیر صفر شده که در *relu* این طور نیست و مقادیر زیر صفر از بین میروند. همچنین به دلیل *smooth* بودن وابستگی شبکه به *learning rate* و وزن های اولیه کم شده و سریع تر میتوان به نقطه بهینه رسید.

ث) میتوان با کنترل بتا خروجی تابع را کنترل کرد به طوری که با زیاد کردن این مقدار میتوان تابع را به *relu* و با کم کردن آن میتوان به تابع خطی میل داد.

ج) با وجود Δ گرادین تابع *mish* نسبت به *smooth swish* تر شده و بهینه سازی بهتر و سریع تر انجام میشود.

۲- الف) اگر مقادیر اولیه به صورت رندوم باشد به طور میانگین خروجی برابر با ۰.۵ میشود که با جایگذاری آن در فرمول دو تابع ضرر مد نظر مقادیر ۰.۷ و ۰.۲۵ بدست میاید و نتیجه دور از انتظاری نیست. ب) خروجی تابع ضرر *mse* عددی بین صفر و یک است ولی در *bce* مقادیر بیشتر هستند، در نتیجه مقادیر میتوانند اختلافات بیشتری داشته باشند. ت) در *bce* همانطور که مشاهده میکنیم، از حدود *epoch* ۶۰ به بعد تابع به سمت اورفیت شدن میرود و مقادیر برای دیتای ولیدیشن افزایش میابد، در نتیجه بهتر است آموزش را در همان لحظات متوقف کنیم، ولی در *mse* همانطور که مشاهده میشود تا *epoch* ها آخر همچنان شیب تابع منفی است و در حال کمن شدن است و حتی میتوان چند *epoch* دیگر آن را جلو برد و نتایج بهتری بدست آورد.

۳- در مدل با آلفای برابر با یک دقت به شدت پایین و چیزی بین ۰.۱۱ تا ۰.۱۷ است که اصلا خوب نیست. علت این موضوع این است که وقتی آلفا یک باشد تابع فعال ساز به صورت خطی خواهد شد و درواقع انگار شبکه ای با یک لایه داریم که طبیعتا دقت پایینی خواهد داشت. همچنین بهترین عملکرد مربوط به آلفای منفی یک است، زیرا هم دقت آن هم در ولیدیشن و هم تست بالا ۹۸ درصد اصن و عملکرد خیلی خوبی داشته است.