

به نام خالق رنگین کمان

ستاره باباجانی 99521109- تمرین سری 1

Subject: ————— « به نام خالق رنگین کمان »

سوال (الف): برای محاسبه احتمال علی ضعیف باشد، به ناسی و قوی ناسی داشته باشد.

ناسی قوی $C=1$ و ناسی ضعیف $C=0$

$$\Rightarrow P(C=0|D) = \frac{P(C=0) \times P(D|C=0)}{P(D)} \quad , \quad P(C=1|D) = \frac{P(C=1) \times P(D|C=1)}{P(D)}$$

در اینجا ما به دنبال محاسبه پستی در محاسبه ناسی قوی و ضعیف هستیم.

$$P(C=0) = P(C=1) = \frac{1}{2} = \frac{1}{2}$$

$$\Rightarrow P(D|C=0) = P(\text{ناسی قوی} | 0) \times P(\text{ناسی ضعیف} | 0) \times P(\text{علی} | 0) \times P(\text{انسانی} | 0)$$

$$= \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} \rightarrow P(C=0|D) = \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} \times \frac{1}{2} = \frac{1}{128} \quad (4\% \text{ صد} 128)$$

$$\Rightarrow P(C=1|D) = \frac{1}{4} \times P(D|C=1) = \frac{1}{4} \times P(\text{ناسی قوی} | 1) \times P(\text{ناسی ضعیف} | 1) \times P(\text{علی} | 1) \times P(\text{انسانی} | 1)$$

$$= \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} = \frac{1}{128} \quad (4\% \text{ صد} 128)$$

در اینجا ما به دنبال محاسبه پستی در محاسبه ناسی قوی و ضعیف هستیم.

$$\Rightarrow P(C=0|D) = \frac{1}{2} \times P(D|C=0) = \frac{1}{2} \times P(C=1|D) = \frac{1}{2} \times P(D|C=1) = \frac{1}{2} \times \frac{1}{128} = \frac{1}{256}$$

این محاسبه با naive Bayes متفاوت است. (این محاسبه naive Bayes است.)

این محاسبه با naive Bayes متفاوت است. (این محاسبه naive Bayes است.)

$$P(D_i|C) = \frac{\text{Freq} + 1}{N + K}$$

number of features in data $\leftarrow N + K$

$$P(D_i|C) = \frac{\text{Freq}(D_i|C) + 1}{N_{\text{class}} + \sqrt{K}}$$

نویس دست دوم فراموش ناسی (d=1):

$$\Rightarrow P(D|C=0) = \frac{3}{14} \times \frac{3}{14} \times \frac{3}{14} \times \frac{1}{14} \times \frac{1}{14} \Rightarrow P(C=0|D) = \frac{1}{4} \times \frac{3 \times 3 \times 3}{14^5} = \frac{1}{4} \times \frac{27}{14^5} \quad (4\% \text{ صد} 14)$$

$$\Rightarrow P(D|C=1) = \frac{2}{13} \times \frac{2}{13} \times \frac{2}{13} \times \frac{1}{13} \times \frac{1}{13} \Rightarrow P(C=1|D) = \frac{1}{4} \times \frac{8}{13^5} = \frac{1}{4} \times \frac{8}{13^5} \quad (4\% \text{ صد} 13)$$

IDEA

در اینجا ما به دنبال محاسبه پستی در محاسبه ناسی قوی و ضعیف هستیم.

مرجع:

<https://towardsdatascience.com/laplace-smoothing-in-na%C3%AFve-bayes-algorithm-9c237a8bdece>

سوال 2: -

سوال ۳: داریم:

$$\Phi(a) = \int_{-\infty}^a N(\theta|0,1) d\theta, \quad N(\theta|0,1) \text{ توزیع نرمال استاندارد}$$

برای محاسبه احتمال فردی y (که می‌تواند ۱ یا ۰ باشد) به شرط ورود \vec{x} و وزن \vec{w} و بایاس b ، خواهیم داشت:

$$P(y|\vec{x}, \vec{w}, b) = \begin{cases} \Phi(\vec{x} \cdot \vec{w} + b) & \text{if } y=1 \\ 1 - \Phi(\vec{x} \cdot \vec{w} + b) & \text{if } y=0 \end{cases}$$

که در اینجا، این مدل برای یک نقطه (\vec{x}_i, y_i) با مقیاس تغییر کند: Likelihood

$$P(y_i|\vec{x}_i, \vec{w}, b) = \begin{cases} \Phi(\vec{x}_i \cdot \vec{w} + b) & \text{if } y_i=1 \\ 1 - \Phi(\vec{x}_i \cdot \vec{w} + b) & \text{if } y_i=0 \end{cases}$$

اگر چند اسلاید، خواهیم داشت: $L(\vec{w}, b)$ تابع هزینه

$$P(y_i|\vec{x}_i, \vec{w}, b) = y_i \cdot \Phi(\vec{x}_i \cdot \vec{w} + b) + (1 - y_i) \cdot (1 - \Phi(\vec{x}_i \cdot \vec{w} + b))$$

← برای محاسبه تابع هزینه، باید از شرطی، تابع هزینه را بنویسیم:

$$L(\vec{w}, b) = -\log(P(y_i|\vec{x}_i, \vec{w}, b))$$

$$L(\vec{w}, b) = \begin{cases} -\log \Phi(\vec{x}_i \cdot \vec{w} + b) & \text{if } y_i=1 \\ -\log(1 - \Phi(\vec{x}_i \cdot \vec{w} + b)) & \text{if } y_i=0 \end{cases}$$

$$L_i(\vec{w}, b) = y_i \cdot \log \Phi(\vec{x}_i \cdot \vec{w} + b) + (1 - y_i) \cdot \log(1 - \Phi(\vec{x}_i \cdot \vec{w} + b))$$

← حالا برای محاسبه هزینه، باید از معادله هزینه را حساب می‌کنیم:

$$L(\vec{w}, b) = \frac{1}{n} \sum_{i=1}^n L_i(\vec{w}, b) = \dots$$

$$L(\vec{w}, b) = -\frac{1}{n} \sum_{i=1}^n y_i \cdot \log(\Phi(\vec{x}_i \cdot \vec{w} + b)) + (1 - y_i) \cdot \log(1 - \Phi(\vec{x}_i \cdot \vec{w} + b))$$

سوال 4: الف) توابع فعال سازی یک جزء ضروری در پرسپترون های چند لایه (MLP) و دیگر معماری های شبکه عصبی هستند. آنها چندین هدف مهم را انجام می دهند، از جمله:

- معرفی غیر خطی بودن: توابع فعال سازی غیر خطی بودن را به مدل وارد می کند، که به شبکه اجازه می دهد تا روابط پیچیده و غیر خطی را در داده ها تقریب زده و یاد بگیرد. بدون توابع فعال سازی غیر خطی، کل شبکه عصبی مانند یک تبدیل خطی عمل می کند و باعث می شود که قادر به یادگیری و نمایش الگوهای پیچیده نباشد.

- توابع پیچیده یادگیری: توابع فعال سازی غیر خطی شبکه عصبی را قادر می سازد تا طیف وسیعی از توابع را تقریب و یاد بگیرد و آنها را برای کارهای مختلف از جمله تشخیص تصویر، پردازش زبان طبیعی و غیره مناسب می کند. ظرفیت مدل سازی توابع پیچیده برای حل مسائل دنیای واقعی بسیار مهم است.

- لایه های پشته ای: MLP ها از چندین لایه نرون تشکیل شده اند. بدون توابع فعال سازی، این لایه ها به طور موثر در یک تبدیل خطی منفرد فرو می روند و شبکه را قدرتمندتر از یک پرسپترون تک لایه نمی سازند. توابع فعال سازی به شما امکان می دهد چندین لایه را روی

هم قرار دهید و شبکه های عصبی عمیقی ایجاد کنید که قادر به یادگیری ویژگی های سلسله مراتبی هستند.

- نمایش پراکنده: برخی از توابع فعال سازی، مانند واحدهای خطی اصلاح شده (ReLU)، می توانند نمایش های پراکنده را در شبکه تشویق کنند. فعال سازی های پراکنده از نظر محاسباتی کارآمد هستند و می توانند به تمرکز شبکه بر روی مرتبط ترین ویژگی ها در داده ها کمک کنند.
- اجتناب از مشکل گرادیان ناپدید شدن: توابع فعال سازی می توانند به کاهش مشکل گرادیان ناپدید شدن که در طول انتشار پس پخش در شبکه های عمیق رخ می دهد، کمک کنند. برخی از توابع فعال سازی، مانند ReLU، برای ورودی های مثبت اشباع نمی شوند و به شیب ها اجازه می دهند راحت تر جریان پیدا کنند و از ناپدید شدن آنها در طول تمرین جلوگیری می کنند.
- مدیریت ورودی های منفی: توابع فعال سازی، مانند توابع Sigmoid و Tanh، می توانند ورودی ها را به محدوده خروجی مورد نظر نگاشت (به عنوان مثال، [0, 1] برای Sigmoid و [-1, 1] برای Tanh)، که می

تواند برای Tanh مفید باشد. کاربردهای خاص مانند طبقه بندی باینری یا رگرسیون محدود.

- منظم‌سازی: برخی از توابع فعال‌سازی، مانند تابع Dropout، می‌توانند با غیرفعال کردن تصادفی بخشی از نورون‌ها در طول تمرین، به‌عنوان نوعی منظم‌سازی عمل کنند. این به جلوگیری از برازش بیش از حد کمک می‌کند و تعمیم مدل را بهبود می‌بخشد.

علاوه بر این، همان‌طور که در اسلاید درس آمده، توابع فعال‌سازی برای هر نورون مشخص می‌کنند که خروجی بخش خطی در چه شرایطی منجر به فعال شدن نورون باید بشود.

ب) در تئوری، می‌توان از هر تابع فعال‌سازی غیر خطی برای پرسپترون چند لایه (MLP) استفاده کرد. انتخاب تابع فعال‌سازی باید به نیازهای خاص شبکه عصبی و ماهیت مشکل بستگی داشته باشد. در اینجا برخی از ملاحظات برای انتخاب یک تابع فعال‌سازی وجود دارد:

- ReLU (واحد خطی اصلاح شده): ReLU به دلیل سادگی و کارایی در بسیاری از موارد یکی از پرکاربردترین توابع فعال‌سازی است. به خصوص

برای شبکه های عمیق مناسب است. با این حال، ممکن است از مشکل "ReLU در حال مرگ" رنج ببرد که در آن برخی از نورون ها در طول تمرین غیرفعال می شوند.

- Sigmoid و Tanh: توابع Sigmoid و Tanh برای کارهایی مانند طبقه بندی باینری مفید هستند، جایی که می خواهید خروجی را در محدوده خاصی محدود کنید (0 تا 1 برای Sigmoid و -1 تا 1 برای Tanh). آنها همچنین در شبکه های عصبی بازگشتی خاص استفاده می شوند.

- Leaky ReLU: Leaky ReLU مشکل "ReLU در حال مرگ" را با اجازه دادن یک گرادیان کوچک برای ورودی های منفی برطرف می کند. این می تواند ثبات تمرین را بهبود بخشد.

- ReLU پارامتریک PReLU: PReLU (PReLU) گونه ای از Leaky ReLU است که در آن شیب قسمت منفی را می توان در طول آموزش یاد گرفت و انعطاف پذیری بیشتری را فراهم می کند.

- واحد خطی نمایی ELU: (ELU) گزینه دیگری است که می تواند مشکل در حال مرگ ReLU را کاهش دهد و نشان داده شده است که در برخی از سناریوها عملکرد خوبی دارد.

- واحد خطی نمایی مقیاس دار SELU: (SELU) یک تابع فعال سازی خودکار است که می تواند به تثبیت فعال سازی ها و گرادیان ها در شبکه های عمیق تحت شرایط خاص کمک کند.

اما اگر بخواهیم به صورت کلی تر راجع به پارامتر هایی که هنگام توابع غیرخطی باید در نظر داشت صحبت کنیم، خواهیم داشت:

- **Task Requirements:** The nature of your problem and its requirements play a significant role in the choice of activation function. For example, for binary classification tasks, the Sigmoid function is commonly used, while for most other tasks, ReLU variants are popular.
- **Gradient Behavior:** Consider the behavior of the activation function's gradient. Some functions, like ReLU, can suffer from the "dying ReLU" problem,

where gradients can become zero for a large portion of the input space. This can affect training stability and slow convergence.

- **Vanishing and Exploding Gradients:** Certain activation functions, like Sigmoid and Tanh, are more prone to vanishing gradients, which can hinder training in deep networks. On the other hand, some functions can lead to exploding gradients, making training difficult. Choose an activation function that mitigates these issues in your specific architecture.
- **Computation Efficiency:** Some activation functions are computationally more efficient than others. Simple functions like ReLU and its variants are faster to compute compared to more complex functions like Sigmoid or Tanh.
- **Training Stability:** The choice of activation function can affect the stability and speed of training. Functions like Leaky ReLU or ELU may provide better

training stability compared to the standard ReLU, especially in deeper networks.

- **Bias Handling:** Consider how the activation function handles bias. Some functions have a natural bias term (e.g., Sigmoid and Tanh), while others do not (e.g., ReLU). This may affect the initial training dynamics and convergence of your model.
- **Output Range:** If your problem requires specific output ranges (e.g., $[0, 1]$ for probabilities in binary classification), choose an activation function that naturally produces outputs in that range.
- **Expressiveness:** Different activation functions have different levels of expressiveness. Some functions, like Swish, have been shown to provide better modeling capabilities for certain tasks. Experiment with different functions to determine which one captures the data's patterns effectively.

- **Architectural Compatibility:** The choice of activation function should be compatible with the overall architecture of your neural network. For example, if you're using a gated architecture like LSTMs or GRUs in a recurrent neural network, the activation functions in those units are specialized for sequential data.
- **Hyperparameter Tuning:** Be prepared to conduct hyperparameter tuning to optimize the choice of activation function along with other hyperparameters. The ideal activation function may vary depending on the specific model architecture and dataset.
- **Empirical Testing:** Ultimately, it's often best to empirically test different activation functions on your specific problem and data. This involves training models with different activation functions and comparing their performance, training speed, and stability.

- Regularization: Some activation functions, like dropout, can also act as a form of regularization. If you need regularization, consider how the activation function may interact with other regularization techniques in your model.

منبع:

<https://chat.openai.com/>

سوال 5: الف)

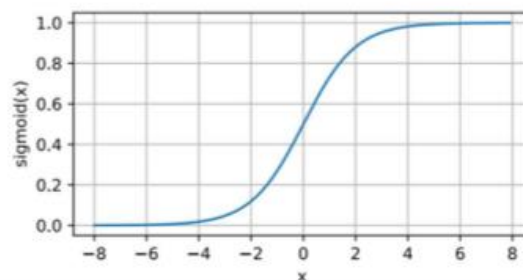
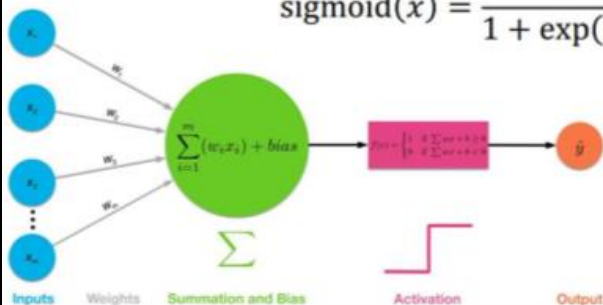
- تابع sigmoid:

Sigmoid

- در اولین شبکه‌های عصبی، دانشمندان علاقه‌مند به مدل‌سازی نورون‌های بیولوژیکی بودند که یا فعال می‌شوند یا نمی‌شوند

- با توسعه روش‌های یادگیری مبتنی بر گرادیان، نیاز بود تا تقریب مشتق‌پذیر استفاده شود

$$\text{sigmoid}(x) = \frac{1}{1 + \exp(-x)}$$

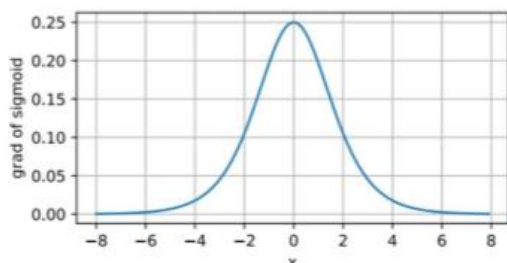




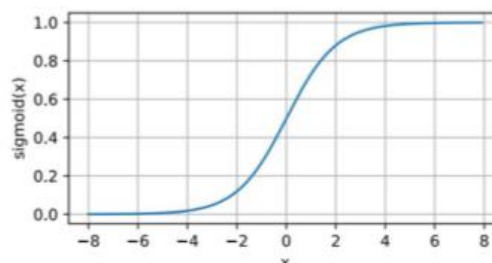
Sigmoid

- در بسیاری از مقادیر، گرادیان نزدیک به صفر است که بهینه‌سازی شبکه‌های عمیق را پیچیده می‌کند
- یکی از ایرادات sigmoid این است که خروجی آن همواره مثبت است

$$\frac{d}{dx} \text{sigmoid}(x) = \text{sigmoid}(x)(1 - \text{sigmoid}(x))$$



$$\text{sigmoid}(x) = \frac{1}{1 + \exp(-x)}$$



Sigmoid یک مقدار واقعی را به عنوان ورودی می‌گیرد و مقدار دیگری را بین 0 و 1 خروجی می‌دهد. تابع فعال سازی sigmoid، ورودی محدوده $(-\infty, \infty)$ را به محدوده (0,1) ترجمه می‌کند.

مزایا:

1. خروجی در محدوده [0, 1] که برای طبقه بندی باینری مناسب است، جایی که می‌توان آن را به عنوان یک احتمال تفسیر کرد.
2. صاف و قابل تمایز است که باعث می‌شود در طول انتشار پس از انتشار به خوبی رفتار کند.
3. به طور تاریخی در شبکه های عصبی و رگرسیون لجستیک استفاده می‌شود.

معایب:

1. از مشکل شیب ناپدید شدن رنج می برد، به خصوص برای ورودی های بسیار مثبت یا بسیار منفی، که می تواند سرعت تمرین را کاهش دهد.
2. خروجی ها صفر محور نیستند، که می تواند بهینه سازی را در شبکه های عمیق چالش برانگیز کند.
3. محدود به وظایف طبقه بندی باینری.

• تابع softmax:

فعال سازی Softmax استفاده کنیم که تعمیم تابع Sigmoid است

$$\mathbf{o} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

$$\hat{y}_i = \text{softmax}(\mathbf{o})_i = \frac{\exp(o_i)}{\sum_{k=1}^q \exp(o_k)}$$

• تابع ضرر متناسب با این تابع فعال سازی، categorical cross-entropy است

$$l(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^q y_i \log \hat{y}_i$$

مشتق softmax

$$\begin{aligned}
 l(\mathbf{y}, \hat{\mathbf{y}}) &= - \sum_{i=1}^q y_i \log \frac{\exp(o_i)}{\sum_{k=1}^q \exp(o_k)} = \sum_{i=1}^q y_i \log \sum_{k=1}^q \exp(o_k) - \sum_{i=1}^q y_i o_i \\
 &= \log \sum_{k=1}^q \exp(o_k) \sum_{i=1}^q y_i - \sum_{i=1}^q y_i o_i = \log \sum_{k=1}^q \exp(o_k) - \sum_{i=1}^q y_i o_i \\
 \partial_{o_i} l(\mathbf{y}, \hat{\mathbf{y}}) &= \frac{\exp(o_i)}{\sum_{k=1}^q \exp(o_k)} - y_i = \text{softmax}(\mathbf{o})_i - y_i
 \end{aligned}$$

تابع Softmax توزیع احتمالات رویداد را بر روی 'n' رویدادهای مختلف محاسبه می کند. به طور کلی، این تابع احتمالات هر کلاس هدف را بر روی تمام کلاس های هدف ممکن محاسبه می کند. بعداً احتمالات محاسبه شده به تعیین کلاس هدف برای ورودی های داده شده کمک می کند.

مزایا:

1. بردار اعداد واقعی را به توزیع احتمال در چندین کلاس تبدیل می کند و آن را برای کارهای طبقه بندی چند کلاسه مناسب می کند.
2. مجموع احتمالات کلاس برابر با 1 است.
3. در لایه خروجی شبکه های عصبی برای طبقه بندی چند کلاسه استفاده می شود.

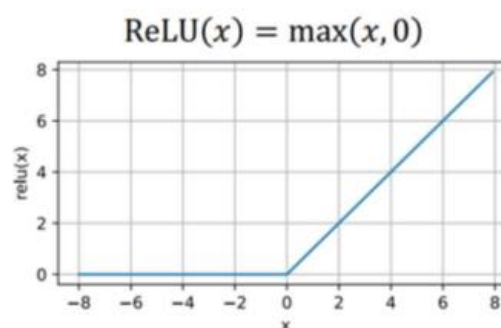
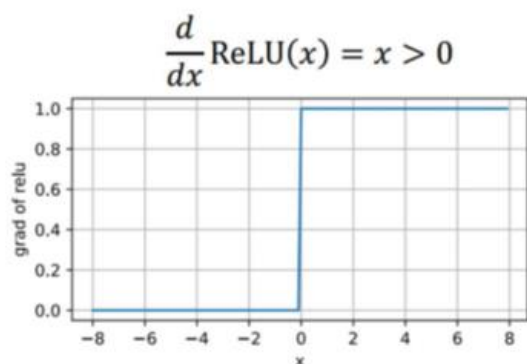
معایب:

1. محاسبه نمایی در تابع softmax می‌تواند منجر به بی‌ثباتی عددی در هنگام برخورد با مقادیر ورودی بسیار بزرگ یا کوچک شود، که به طور بالقوه باعث مشکلات سرریز یا زیر جریان می‌شود.
2. خروجی softmax بسیار به مقادیر ورودی وابسته است و آن را به مقادیر پرت حساس می‌کند.

• تابع ReLU:

Rectified Linear Unit (ReLU)

- یک تابع غیرخطی بسیار ساده و پرکاربرد است
- بهینه‌سازی ReLU بسیار آسان است زیرا بسیار شبیه به واحدهای خطی است
- مشتق ReLU برای مقادیری که فعال است همواره بزرگ است



فرمول ساده است: $\max(0, z)$. علی‌رغم نامش، Rectified Linear Units، خطی نیست و همان مزایای Sigmoid را دارد اما عملکرد بهتری دارد.

مزایا:

1. از نظر محاسباتی کارآمد است، زیرا شامل آستانه گذاری ساده است.
2. از مشکل ناپدید شدن گرادیان برای ورودی های مثبت رنج نمی برد، که می تواند آموزش در شبکه های عمیق را سرعت بخشد.
3. پراکندگی در فعال سازی را ترویج می کند و باعث می شود شبکه بر روی ویژگی های مرتبط تمرکز کند.
4. به طور گسترده مورد استفاده قرار می گیرد و عملکرد بسیار خوبی در کارهای مختلف به ویژه یادگیری عمیق نشان داده است.

معایب:

1. ممکن است از مشکل "ReLU در حال مرگ" رنج ببرد، جایی که برخی از نورون ها می توانند غیرفعال شوند (همیشه خروجی صفر دارند) و هرگز در طول تمرین برای ورودی های منفی بهبود نمی یابند.
2. خروجی ها محدود نیستند، که می تواند در موارد خاصی که خروجی های محدود مورد نیاز است، مشکلاتی ایجاد کند.
3. صفر محور نیست، که می تواند منجر به مشکلات همگرایی در برخی از معماری های شبکه شود.

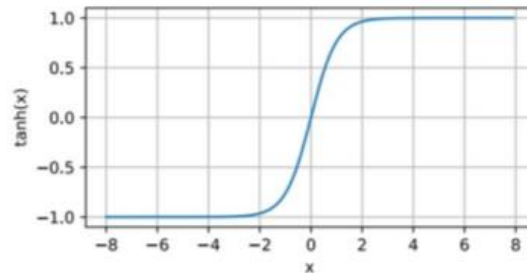
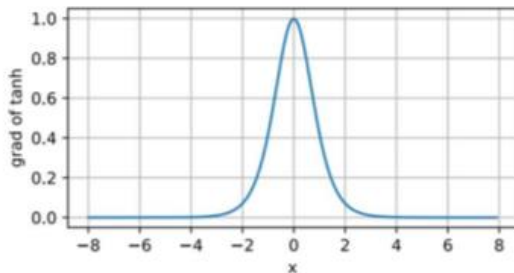
• تابع Tanh:

Tanh

- مشابه با تابع Sigmoid خروجی آن محدود است اما به بازه -1 تا $+1$
- اشباع شدن تابع \tanh ، بهینه‌سازی را دشوار می‌کند و از لحاظ محاسباتی هم به دلیل \exp پرهزینه است

$$\frac{d}{dx} \tanh(x) = 1 - \tanh^2(x)$$

$$\tanh(x) = \frac{1 - \exp(-2x)}{1 + \exp(-2x)} = 2\sigma(2x) - 1$$



تابع \tanh یکی دیگر از تابع‌های ممکن است که می‌تواند به عنوان یک تابع فعال‌سازی غیر خطی بین لایه‌های یک شبکه عصبی استفاده شود. چند چیز مشترک با عملکرد فعال‌سازی سیگموئید دارد. برخلاف تابع سیگموئید که مقادیر ورودی را بین 0 و 1 ترسیم می‌کند، \tanh مقادیر بین -1 و 1 را ترسیم می‌کند. مشابه تابع سیگموئید، یکی از ویژگی‌های جالب تابع \tanh این است که مشتق \tanh را می‌توان در بیان کرد. مزایا:

1. خروجی در محدوده $[-1, 1]$ ، که به کاهش مشکل گرادیان ناپدید شدن در مقایسه با Sigmoid کمک می‌کند.
2. صفر محور، آن را برای بهینه‌سازی و همگرایی در شبکه‌های عمیق مناسب‌تر می‌کند.

3. صاف و قابل تمایز، تسهیل انتشار پس زمینه.

معایب:

1. هنوز مستعد مشکل گرادیان ناپدید شدن است، به ویژه برای ورودی

های بسیار مثبت یا بسیار منفی.

2. عملیات نمایی آن را از نظر محاسباتی گرانتر از ReLU می کند.

3. خروجی ها ممکن است همیشه برای همه انواع وظایف مناسب

نباشند، زیرا محدود هستند.

حال به مقایسه این توابع میپردازیم:

Activation Function	Output Range	Advantages	Disadvantages
Sigmoid	[0, 1]	- Suitable for binary classification. - Smooth and differentiable. - Historically used in neural networks.	- Suffers from vanishing gradient, especially for extreme inputs. - Outputs not zero-centered. - Limited to binary classification.
Softmax	[0, 1] (sums to 1)	- Converts input into a probability distribution for multi-class classification. - Enforces sum of probabilities to be 1.	- Exponential computation can lead to numerical instability. - Sensitive to outliers in input values.
ReLU (Rectified Linear Unit)	[0, ∞)	- Computationally efficient and fast due to simple thresholding. - Mitigates vanishing gradient problem for positive inputs. - Promotes sparse activations. - Widely used with good performance in deep learning.	- "Dying ReLU" problem for negative inputs (some neurons become inactive). - Outputs are unbounded. - Not zero-centered.
Tanh (Hyperbolic Tangent)	[-1, 1]	- Outputs in [-1, 1], reducing vanishing gradient issues compared to Sigmoid. - Zero-centered, aiding convergence. - Smooth and differentiable.	- Still susceptible to vanishing gradient, especially for extreme inputs. - More computationally expensive compared to ReLU. - Outputs may not always be suitable for all tasks.

منابع:

<https://chat.openai.com/>

<https://www.analyticsvidhya.com/blog/2021/04/activation-functions-and-their-derivatives-a-quick-complete-guide>

ب) حال در این بخش میخواهیم این توابع را پیاده سازی کنیم. در 3 سلول اول، کتابخانه ها مورد نیاز ایمپورت شد و دیتا ساخته شد و یک تابع برای نمایش توابع طراحی شد:

```
[1] 1 # import libraries
      2 import numpy as np
      3 import torch
      4 import torch.nn as nn
      5 import matplotlib.pyplot as plt

[2] 1 # create some dummy data
      2 test_data_np = np.linspace(-10.0, 10.0, num=500)
      3 test_data_torch = torch.tensor(test_data_np, dtype=torch.float64)

[3] 1 # define a function for plotting
      2 def plot_function(x, y):
      3     plt.plot(x, y)
```

این تابع نیز برای تست توابع نوشته شده استفاده میشود:

```
▶ 1 # define a function to test the two activation functions
    2 def test(x, function1, function2):
    3     output_fn1 = function1(x)
    4     output_fn2 = function2(x)
    5     return torch.allclose(output_fn1, output_fn2), output_fn1
```

حال به بررسی کدهای جدید زده شده میپردازیم:

1. تابع sigmoid: طبق فرمول داده شده و با استفاده از تابع exp داریم:

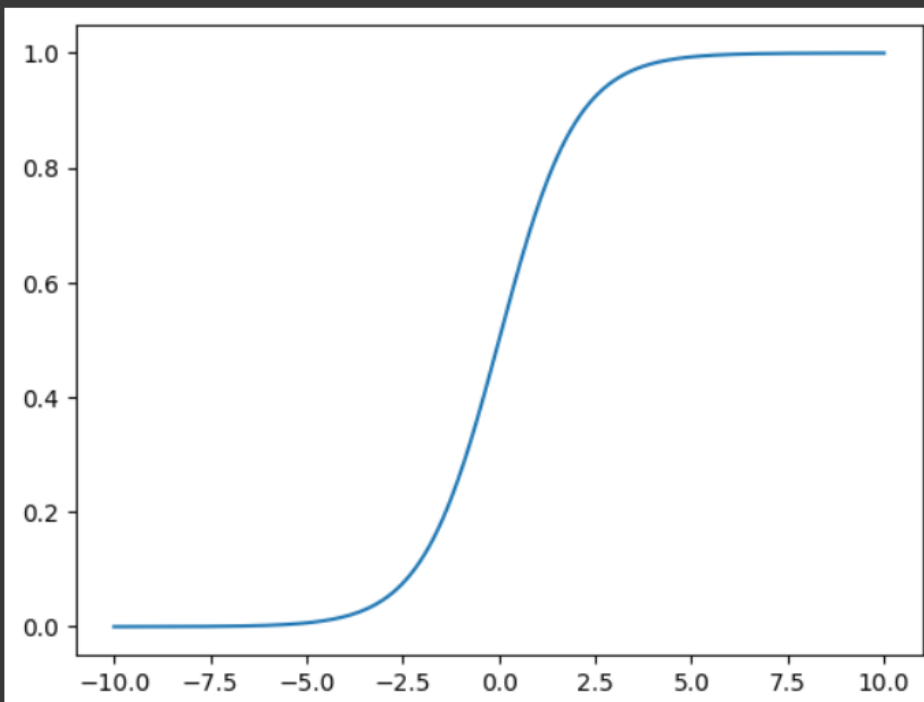
```
[6] 1 # implement the sigmoid activation function
    2 ##### TO DO #####
    3 def my_sigmoid(x):
    4     return 1 / (1 + torch.exp(-x))
```

حال خروجی تست این تابع و سپس شکل آن را میبینیم:

```
1 # test your function against the sigmoid activation function from torch
2 result_sigmoid, output_sigmoid = test(test_data_torch, my_sigmoid, nn.Sigmoid())
3 print(result_sigmoid)
```

True

```
1 # plot the sigmoid activation function
2 plot_function(test_data_np, output_sigmoid)
```



2. تابع softmax: طبق فرمول و با استفاده از توابع آماده exp, sum داریم: (تابع sum جمع را حساب میکند)

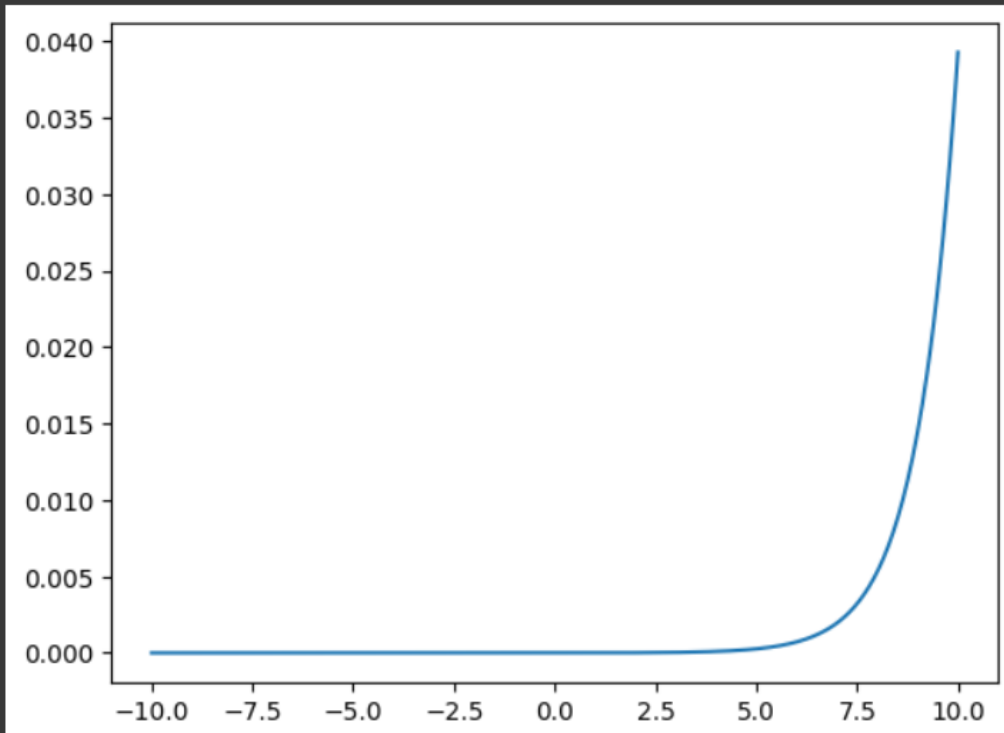
```
1 # implement the softmax activation function
2 ##### TO DO #####
3 def my_softmax(x):
4     return torch.exp(x) / torch.sum(torch.exp(x))
```

حال خروجی تست این تابع و سپس شکل آن را میبینیم:

```
9] 1 # test your function against the softmax activation function from torch
    2 result_softmax, output_softmax = test(test_data_torch, my_softmax, nn.Softmax())
    3 print(result_softmax)

True
```

```
1 # plot the softmax activation function
2 plot_function(test_data_np, output_softmax)
```



3. تابع ReLU: با استفاده از فرمول و تابع آماده max (برای ماکسیمم گرفتن) داریم: (تابع tensor مقدار تنسور 0 است)

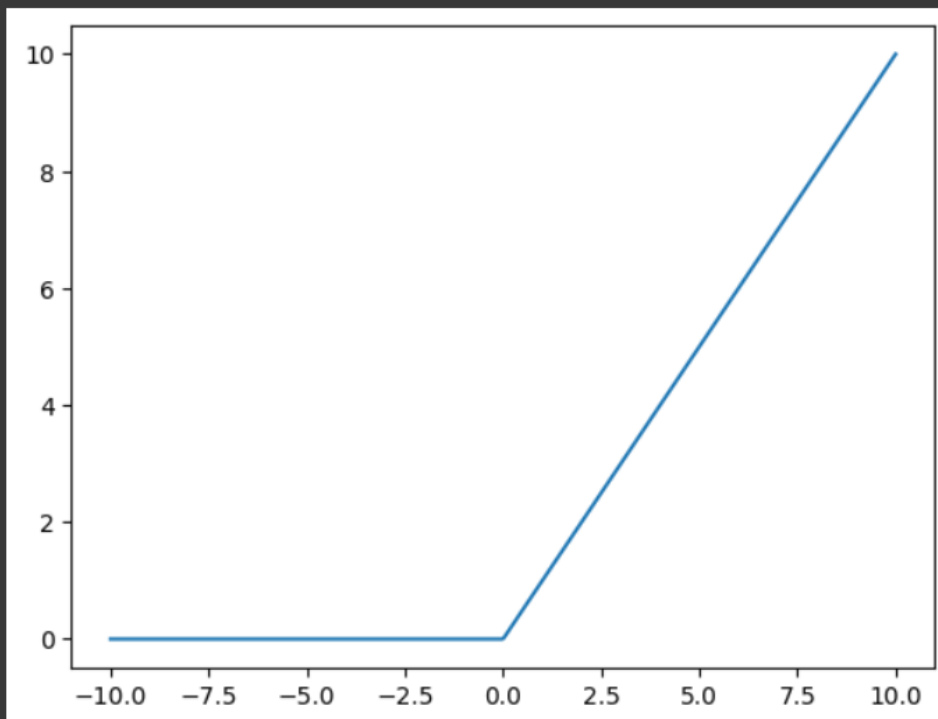
```
1 # implement the ReLU activation function
2 ##### TO DO #####
3 def my_relu(x):
4     return torch.max(torch.tensor(0), x)
```


حال خروجی تست این تابع و سپس شکل آن را میبینیم:

```
1 # test your function against the ReLU activation function from torch
2 result_relu, output_relu = test(test_data_torch, my_relu, nn.ReLU())
3 print(result_relu)
```

True

```
1 # plot the ReLU activation function
2 plot_function(test_data_np, output_relu)
```



4. تابع Leaky ReLU: طبق فرمول و با در نظر گرفتن مقدار 0.05 برای

آلفا (این مقدار قابل تغییر است) خواهیم داشت:

```

1 # implement the Leaky ReLU activation function
2 ##### TO DO #####
3 def my_leaky_relu(x, alpha=0.05): # adjust the value of alpha
4     return torch.max(alpha * x, x)

```

حال خروجی تست این تابع و سپس شکل آن را میبینیم:

```

1 # test your function against the Leaky ReLU activation function from torch
2 result_leaky_relu, output_leaky_relu = test(test_data_torch, my_leaky_relu, nn.LeakyReLU())
3 print(result_leaky_relu)

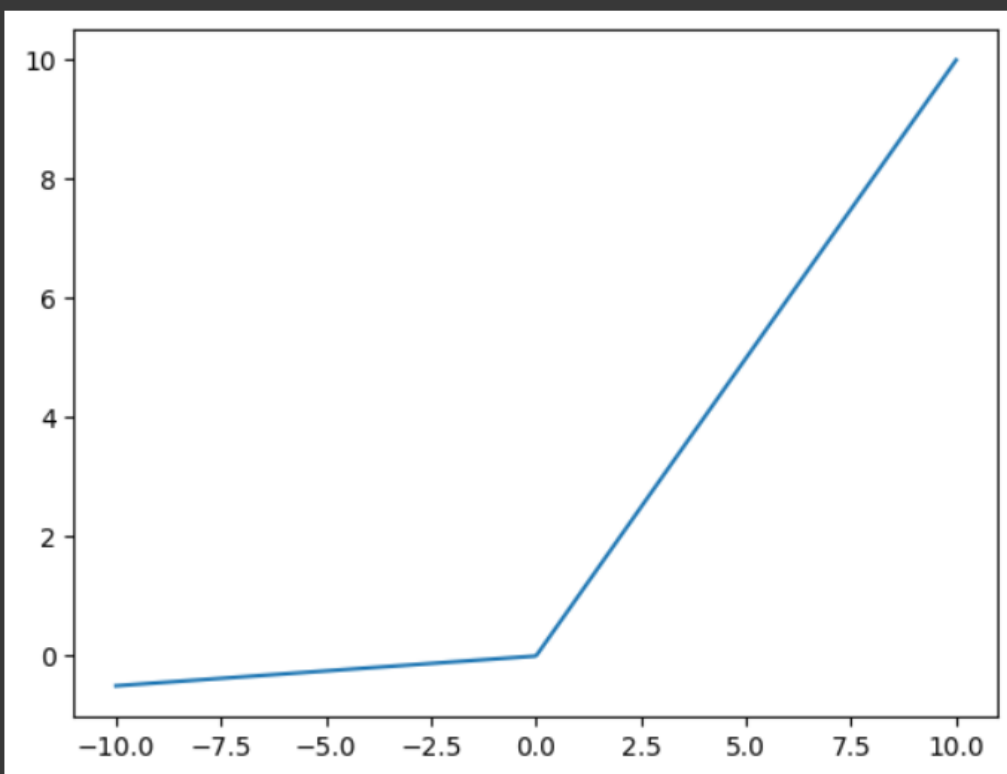
```

False

```

1 # plot the Leaky ReLU activation function
2 plot_function(test_data_np, output_leaky_relu)

```



5. تابع Tanh: طبق فرمول و استفاده از تابع exp خواهیم داشت:

```

1 # implement the tanh activation function
2 ##### TO DO #####
3 def my_tanh(x):
4     return (1 - torch.exp(-2 * x)) / (1 + torch.exp(-2 * x))

```

حال خروجی تست این تابع و سپس شکل آن را میبینیم:

```

1 # test your function against the tanh activation function from torch
2 result_tanh, output_tanh = test(test_data_torch, my_tanh, nn.Tanh())
3 print(result_tanh)

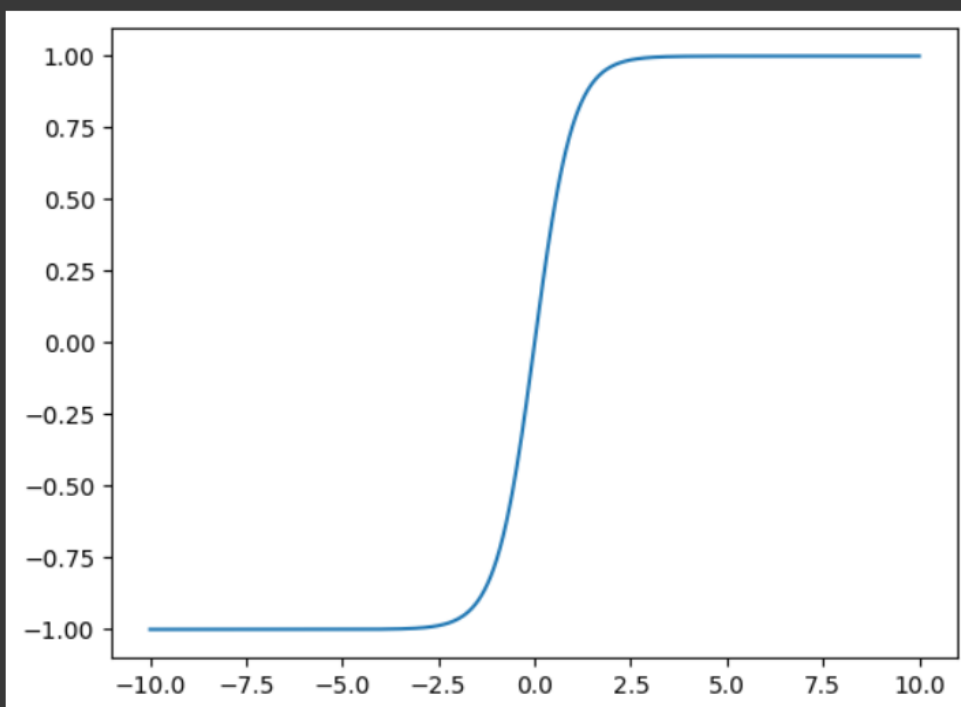
```

True

```

1 # plot the tanh activation function
2 plot_function(test_data_np, output_tanh)

```



همان طور که مشاهده میشود، تمامی توابع درست پیاده سازی شده اند.

ج) معماری که اینجانب معرفی میکنم به این شرح است:

1. تعداد لایه ها و دلیل:

- لایه ورودی: 64 نورون (یکی برای هر پیکسل). لایه ورودی باید به اندازه تعداد پیکسل در تصاویر نورون داشته باشد.
- لایه پنهان: 128 نورون (یا هر مقدار معقول دیگر). یک لایه پنهان برای این کار ساده کافی است. 128 نورون به عنوان یک نقطه شروع معقول انتخاب می شوند و می توان آن را بر اساس عملکرد تنظیم کرد.
- لایه خروجی: 3 نورون (یکی برای هر کلاس). از آنجایی که سه کاراکتر مختلف برای طبقه بندی داریم، لایه خروجی باید سه نورون داشته باشد، یکی برای هر کلاس.

2. تعداد نورون ها و دلیل:

تعداد نورون ها در لایه های ورودی و خروجی بر اساس ماهیت داده ها و وظیفه طبقه بندی تعیین می شود.

تعداد نورون ها در لایه پنهان تا حدودی دلخواه است و می توان آن را در طول آزمایش تنظیم کرد. 128 نورون به عنوان نقطه شروع انتخاب می شوند تا به شبکه اجازه دهند الگوهای نسبتاً پیچیده ای را بیاموزند.

3. تابع و دلیل فعال سازی:

از فعال سازی ReLU (واحد خطی اصلاح شده) برای لایه پنهان استفاده میکنیم. ReLU یک انتخاب رایج است زیرا به جریان گرادیان کمک می کند و می تواند غیرخطی ها را به خوبی مدیریت کند.

از یک تابع فعال سازی Softmax برای لایه خروجی برای تولید احتمالات برای هر کلاس استفاده میکنیم. Softmax معمولا برای مسائل طبقه بندی چند کلاسه استفاده می شود.

4. تابع ضرر و علت آن:

از Cross-Entropy Loss برای مسائل طبقه بندی با کلاس های متعدد استفاده میکنیم. از دست دادن متقابل آنتروپی یک انتخاب مناسب برای طبقه بندی چند کلاسه است و شبکه را تشویق می کند تا احتمالاتی نزدیک به 0 یا 1 برای کلاس صحیح تولید کند.

(د) حال می خواهیم کد معماری تعریف شده را بنویسیم.

1. ابتدا تمامی کتابخانه هایی که به آنها نیاز خواهیم داشت را ایمپورت میکنیم:

```

1 import os
2 from PIL import Image
3 import torch
4 import torch.nn as nn
5 import torch.optim as optim
6 import torchvision.transforms as transforms
7 import matplotlib.pyplot as plt

```

2. حال مدلی که ذکر کردیم را تعریف میکنیم: این مدل شامل یک لایه پنهان و یک لایه خروجی است که ورودی و خروجی آنها در شکل قابل رویت است. ورودی لایه میانی تعداد پیکسل ها یعنی 256 و خروجی آن که ورودی لایه بعدی است، 128 است. لایه پایانی نیز خروجی 3 به علت سه کلاس بودن مسئله دارد.

تابع forward نیز forward passing را مشخص میکند و تابع فعال سازی هر لایه را تعریف میکند.

```

1 # Define the MLP model
2 # The super() function should be called to initialize the parent class, and it should not be used to call the __init__ method.
3 class SimpleMLP(nn.Module):
4     def __init__(self):
5         super(SimpleMLP, self).__init__()
6         self.fc1 = nn.Linear(256, 128) # Hidden Layer with 128 neurons
7         self.fc2 = nn.Linear(128, 3) # Output Layer with 3 neurons (3 classes)
8
9     def forward(self, x):
10        x = torch.relu(self.fc1(x)) # ReLU activation for the hidden layer
11        x = torch.softmax(self.fc2(x), dim=1) # Softmax activation for the output layer
12        return x

```

3. حال باید تصویر ها را به نحوی reshape کرد که برای مدل و تنسور قابل فهم باشد. سپس دیتای آموزش و ارزیابی تعریف میشوند.(کامنت هر خط در شکل قابل رویت است):

```
1 # with the help of chat GPT:
2 def load_image(image_file):
3     img = Image.open(image_file)
4     # Resize the image to 8x8
5     img = img.resize((8, 8))
6     # Convert the image to a tensor and flatten it into a 1D array
7     img_tensor = transforms.ToTensor()(img).view(256)
8     return img_tensor
9
10 image_files = ["Q5_1.png", "Q5_2.png", "Q5_3.png"]
11
12 # Load and preprocess the images
13 image_data = [load_image(image_file) for image_file in image_files]
14 labels = [0, 1, 2]
15
16 train_data = image_data
17 train_labels = labels
18 test_data = image_data
19 test_labels = labels
20
```

لیبل های کلاس به ترتیب ورودی ها 1و2و3 گذاشته شد و چون دیتای کمی داریم، دیتای آموزش و ارزیابی یکی شد که همان طور که انتظار میرود accuracy برای داده test، 100 خواهد شد.

4. حال مدل را صدا زده و training را آغاز میکنیم:

```
1 model = SimpleMLP()
2
3 criterion = nn.CrossEntropyLoss() # Cross-Entropy Loss for multi-class classification
4 optimizer = optim.SGD(model.parameters(), lr=0.01) # Stochastic Gradient Descent optimizer
5
6 #these numbers are custom
7 num_epochs = 100
8 batch_size = 3 # Since we have three images
9 train_loss_values = []
10
11 #Training
12 for epoch in range(num_epochs):
13     epoch_loss = 0.0 # Initialize loss for this epoch
14     for i in range(0, len(train_data), batch_size):
15         inputs = torch.stack(train_data[i:i+batch_size])
16         labels = torch.tensor(train_labels[i:i+batch_size])
17
18         optimizer.zero_grad() # Zero the gradients
19         outputs = model(inputs) # Forward pass
20         loss = criterion(outputs, labels) # Calculate the loss
21
22         loss.backward() # Backpropagation for optimizing
23         optimizer.step() # Update the model's parameters
24         epoch_loss += loss.item()
25
26 # Calculate the average loss for this epoch
27 avg_epoch_loss = epoch_loss / len(train_data)
28 train_loss_values.append(avg_epoch_loss) # Append to the list
29 print(f"Epoch [{epoch+1}/{num_epochs}] - Loss: {loss.item()}")
30
31
```

همان طور که مشاهده میشود ابتدا مدل صدا زده شد و سپس تابع ضرر و تابع optimizer تعریف شد. تعداد ایپوک ها بصورت دلخواه 100 گذاشته شد و سائز هر batch هم به دلیل سه عکس، 3 گذاشته شد.

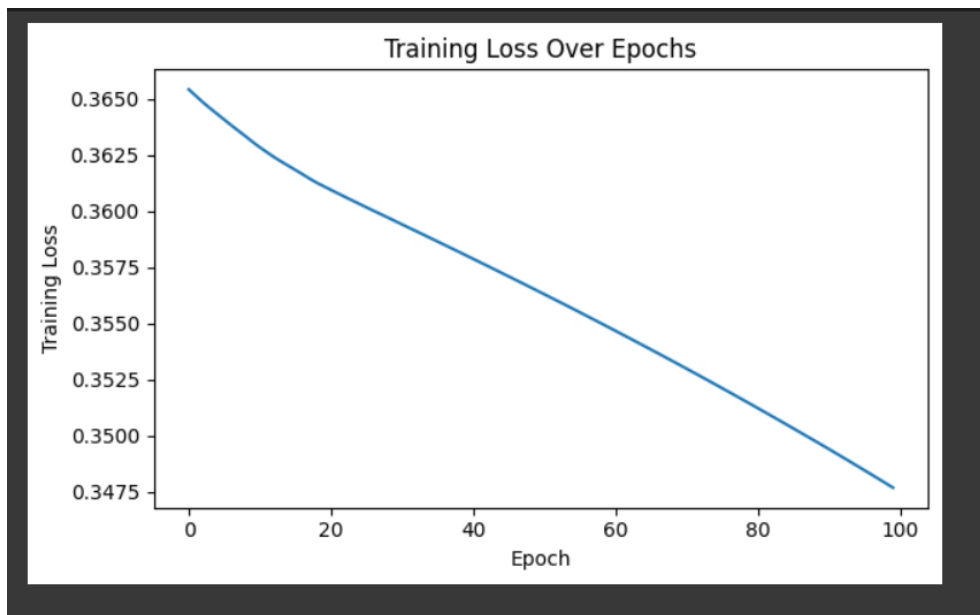
در کد طبق تعداد ایپوک ها حلقه داریم و سپس هر عکس به مدل داده میشود و فرآیند forward passing انجام میشود. در نهایت برای اپدیت کردن وزن ها، از optimizer استفاده میشود. (از back propagation استفاده شده

است) سپس برای اینکه خروجی بهتری داشته باشیم در هر ایپوک مقدار `loss` چاپ شد:

```
Epoch [1/100] - Loss: 1.10541570186615
Epoch [2/100] - Loss: 1.1038767099380493
Epoch [3/100] - Loss: 1.1025326251983643
Epoch [4/100] - Loss: 1.101553201675415
Epoch [5/100] - Loss: 1.100767970085144
Epoch [6/100] - Loss: 1.100049614906311
Epoch [7/100] - Loss: 1.0993293523788452
Epoch [8/100] - Loss: 1.0986448526382446
Epoch [9/100] - Loss: 1.0979958772659302
Epoch [10/100] - Loss: 1.0974016189575195
Epoch [11/100] - Loss: 1.0968527793884277
Epoch [12/100] - Loss: 1.0962334871292114
Epoch [13/100] - Loss: 1.095621109008789
Epoch [14/100] - Loss: 1.095030426979065
Epoch [15/100] - Loss: 1.0944331884384155
Epoch [16/100] - Loss: 1.093700647354126
Epoch [17/100] - Loss: 1.092957615852356
Epoch [18/100] - Loss: 1.0922226905822754
Epoch [19/100] - Loss: 1.0915024280548096
Epoch [20/100] - Loss: 1.090773105621338
Epoch [21/100] - Loss: 1.0900579690933228
Epoch [22/100] - Loss: 1.0893398523330688
```

برای اینکه درک بهتری از خروجی داشته باشیم، از کتابخانه `matplotlib` استفاده کردم. کد زده شده و خروجی آن به شرح زیر است:

```
1 # Plot the training loss
2 plt.figure(figsize=(12, 4))
3 plt.subplot(1, 2, 1)
4 plt.plot(train_loss_values)
5 plt.xlabel("Epoch")
6 plt.ylabel("Training Loss")
7 plt.title("Training Loss Over Epochs")
8
9 plt.tight_layout()
10 plt.show()
```



5. در نهایت مدل آموزش دیده را ارزیابی میکنیم. چون داده آموزش و تست یکی است، accuracy نهایی 100 میشود:

```
31
32 # Testing
33 correct = 0
34 total = 0
35 with torch.no_grad():
36     for i in range(len(test_data)):
37         inputs = test_data[i].unsqueeze(0) # Convert a single image to a 1-sample batch
38         labels = test_labels[i]
39
40         outputs = model(inputs)
41         _, predicted = torch.max(outputs, 1)
42         total += 1
43         if predicted == labels:
44             correct += 1
45
46 print("-----")
47 print(f"Accuracy on test data: {100 * correct / total}%")
```

```
-----
Accuracy on test data: 100.0%
```

سوال 6: ترکیب آستانه 0.5 با تابع فعال سازی ReLU و تابع فعال سازی سیگموئید نهایی در یک شبکه عصبی برای طبقه بندی باینری می تواند برخی از چالش ها و مشکلات بالقوه را معرفی کند:

1. غیر یکنواختی ReLU:

تابع فعال سازی ReLU (واحد خطی اصلاح شده) در 0 قابل تمایز نیست. هنگام استفاده از ReLU، ممکن است تغییرات ناگهانی در خروجی شبکه در حدود آستانه 0.5 ایجاد شود که می تواند فرآیند آموزش را چالش برانگیزتر کند. به طور خاص، مشتق ReLU در 0 تعریف نشده است، و این می تواند بر الگوریتم های بهینه سازی مبتنی بر گرادیان تأثیر بگذارد.

2. چالش های آموزشی:

ترکیب ReLU با یک تابع فعال سازی سیگموئید نهایی، عدم یکنواختی در خروجی شبکه را معرفی می کند. آموزش یک شبکه عصبی با توابع فعال سازی غیر هموار می تواند منجر به مشکلات همگرایی، تمرین کند یا حتی گیر کردن در حداقل های محلی شود. ممکن است برای مدیریت این چالش ها به تکنیک های بهینه سازی پیشرفته تری نیاز باشد.

3. ناپدید شدن و انفجار گرادیان:

استفاده از ReLU به عنوان یک تابع فعال سازی میانی می تواند منجر به مشکلات ناپدید شدن گرادیان یا انفجار در حین انتشار پس انداز شود. تابع فعال سازی سیگموئید نهایی به له کردن خروجی شبکه تا محدوده $[0, 1]$ کمک می کند، اما لایه های میانی ممکن است همچنان با گرادیان ها مشکل داشته باشند که تنظیم اولیه وزن مناسب و نرمال سازی دسته ای می تواند به کاهش این مشکلات کمک کند.

4. مدیریت داده های نامتعادل:

هنگام برخورد با مجموعه داده های نامتعادل، آستانه ثابت 0.5 ممکن است منجر به پیش بینی های مغرضانه شود که به نفع طبقه اکثریت است. مدیریت موثر عدم تعادل کلاس با استفاده از یک آستانه ثابت می تواند چالش برانگیز باشد. تکنیک هایی مانند یادگیری حساس به هزینه یا استفاده از مقادیر آستانه متفاوت برای کلاس های مختلف ممکن است ضروری باشد.

5. تفسیرپذیری مدل:

یک آستانه ثابت می تواند مدل را کمتر قابل تفسیر کند، زیرا ممکن است مشخص نباشد که چرا یک آستانه خاص انتخاب شده است و چگونه بر تصمیمات مدل تأثیر می گذارد. تنظیم آستانه برای موارد استفاده مختلف یا تفسیر خروجی مدل به روشی معنادار می تواند چالش برانگیز باشد.

برای پرداختن به این چالش ها، ممکن است لازم باشد که معماری های شبکه، توابع فعال سازی و مقادیر آستانه ای مختلف را آزمایش کنیم تا بهترین پیکربندی را برای کار طبقه بندی خاص پیدا کنیم. علاوه بر این، برای ارزیابی عملکرد مدل باید از معیارهای ارزیابی استفاده کرد که برای انتخاب آستانه قوی تر هستند، مانند منحنی های فراخوان دقیق یا منحنی های ROC.

سوال 7: الف) مهمترین تفاوت بین یادگیری ماشین و یادگیری عمیق در معماری و عمق شبکه های عصبی مورد استفاده در این رویکردها نهفته است. در اینجا یک تمایز اساسی وجود دارد:

1. یادگیری ماشینی (ML):

- یادگیری ماشینی حوزه وسیع تری است که تکنیک ها و الگوریتم های مختلفی برای تحلیل و مدل سازی داده ها را در بر می گیرد.
- الگوریتم های ML شامل رگرسیون خطی، درخت های تصمیم، ماشین های بردار پشتیبان، k-نزدیک ترین همسایه ها، جنگل های تصادفی و غیره هستند.
- این الگوریتم ها اغلب شامل مهندسی ویژگی می شوند، جایی که دانش خاص دامنه برای انتخاب و تبدیل ویژگی های ورودی مربوطه استفاده می شود.
- مدل های ML معمولاً کم عمق هستند، به این معنی که دارای تعداد کمی لایه یا پارامتر هستند.
- انتخاب ویژگی و مهندسی مراحل بسیار مهمی در یادگیری ماشینی سنتی هستند تا داده ها را برای مدل قابل درک کنند.
- مدل های ML به طور کلی برای طیف گسترده ای از کاربردها، از جمله طبقه بندی، رگرسیون، خوشه بندی و سیستم های توصیه مناسب هستند.

2. یادگیری عمیق (DL):

- یادگیری عمیق زیرمجموعه ای از یادگیری ماشینی است که بر روی شبکه های عصبی با لایه های متعدد (شبکه های عصبی عمیق) تمرکز می کند.
- الگوریتم های DL شامل شبکه های عصبی پیش خور، شبکه های عصبی کانولوشنال (CNN)، شبکه های عصبی بازگشتی (RNN) و غیره هستند.
- مدل های یادگیری عمیق به طور خودکار ویژگی های سلسله مراتبی را از داده های خام یاد می گیرند و نیاز به مهندسی ویژگی های گسترده را از بین می برند.
- مدل های DL با عمقشان مشخص می شوند، با لایه های مخفی متعدد که می توانند الگوهای پیچیده و سلسله مراتبی را در داده ها یاد بگیرند و نشان دهند.
- DL در کارهایی مانند تشخیص تصویر، پردازش زبان طبیعی، تشخیص گفتار و سایر زمینه هایی که داده های خام را می توان مستقیماً پردازش کرد، بسیار موفق بوده است.
- یادگیری عمیق اغلب به حجم زیادی از داده ها و منابع محاسباتی قابل توجهی مانند GPU ها برای آموزش موثر شبکه های عصبی عمیق نیاز دارد.

به طور خلاصه، تفاوت اصلی در عمق شبکه های عصبی و سطح مهندسی ویژگی مورد نیاز است. یادگیری ماشین سنتی بر مهندسی ویژگی های دستی تکیه دارد و معمولاً از مدل های کم عمق تر استفاده می کند. از سوی دیگر، یادگیری عمیق شامل شبکه های عصبی عمیقی است که به طور خودکار ویژگی ها را از داده های خام یاد می گیرند و آن را به ویژه برای کارهایی با مجموعه داده های بزرگ و الگوهای پیچیده، مانند تشخیص تصویر و گفتار، قدرتمند می کنند. با این حال، یادگیری عمیق همچنین به داده ها و منابع محاسباتی قابل توجهی نیاز دارد که می تواند در برخی از برنامه ها محدودیت باشد.

ب) در یک شبکه عصبی عمیق با 16 لایه، معمولاً لایه های نهایی، از جمله لایه یازدهم، برای گرفتن نتیجه نهایی راحت تر هستند. در اینجا دلیل آن است:

1. ناپدید شدن گرادیان: یکی از چالش های آموزش شبکه های عصبی بسیار عمیق، مشکل گرادیان ناپدید شدن است. همانطور که به عمق شبکه می روید، گرادیان ها می توانند بسیار کوچک شوند و به روز رسانی پارامترهای لایه های قبلی را دشوار می کند. این بدان معنی است که لایه های اولیه (آنهايي که نزدیک به ورودی هستند) ممکن است ویژگی های پیچیده و سطح بالا را به طور موثر یاد نگیرند.

2. یادگیری ویژگی سلسله مراتبی: شبکه های عصبی عمیق برای یادگیری ویژگی های سلسله مراتبی طراحی شده اند. در این سلسله مراتب، لایه های اولیه اغلب ویژگی های سطح پایین، مانند لبه ها و الگوهای ساده را به تصویر می کشند، در حالی که لایه های بعدی ویژگی های سطح بالاتر و انتزاعی تری را به تصویر می کشند. لایه های نهایی مسئول ترکیب این ویژگی های سطح بالا برای تصمیم گیری یا پیش بینی نهایی هستند.

3. قدرت بازنمایی: لایه های عمیق تر در یک شبکه عصبی قدرت نمایش بالاتری دارند، زیرا می توانند الگوهای پیچیده و انتزاعی فزاینده ای را به تصویر بکشند. به همین دلیل است که لایه های نهایی اغلب برای گرفتن نتیجه نهایی مناسب هستند، زیرا ظرفیت بیان بیشتری دارند.

4. ویژگی های Task-Specific: در بسیاری از موارد، لایه های نهایی مختص کار هستند. به عنوان مثال، در یک شبکه عصبی کانولوشن (CNN) برای طبقه بندی تصویر، لایه های نهایی از لایه های کاملاً متصل تشکیل شده اند که ویژگی های آموخته شده در لایه های قبلی را برای پیش بینی کلاس ترکیب می کنند.

در یک شبکه عصبی عمیق معمولی، لایه‌های عمیق‌تر، از جمله لایه یازدهم، برای گرفتن اطلاعات سطح بالا و مختص کار و در نتیجه نتیجه نهایی راحت‌تر هستند.

ج) انتخاب بین شبکه‌های عمیق‌تر و گسترده‌تر به مشکل خاص، داده‌های موجود و منابع محاسباتی بستگی دارد. در اینجا یک دستورالعمل کلی وجود دارد:

اگر داده‌های موجود محدود باشد، معمولاً گرایش به شبکه‌های گسترده‌تر کارآمدتر است. شبکه‌های گسترده‌تر می‌توانند انواع ویژگی‌های ساده را به خود اختصاص دهند و در زمانی که داده‌ها کمیاب هستند، کمتر مستعد بیش از حد برازش هستند. آنها در مقایسه با شبکه‌های بسیار عمیق از نظر داده کارآمدتر هستند و به پارامترهای کمتری نیاز دارند.

با این حال، اگر به مقدار قابل توجهی از داده‌ها دسترسی داشته باشیم، به طور کلی انتخاب شبکه‌های عمیق‌تر کارآمدتر است. شبکه‌های عمیق‌تر قادر به یادگیری ویژگی‌های سلسله‌مراتبی و پیچیده‌تر هستند، و آنها را برای کارهایی که شامل الگوها و روابط پیچیده هستند، مناسب می‌سازد. در حالی که آموزش آنها می‌تواند چالش برانگیز باشد، اما می‌توانند تقریب‌های دقیق‌تری را برای توابع با داده‌های فراوان ارائه دهند.

در عمل، ما همچنین می‌توانیم معماری‌های ترکیبی را در نظر بگیریم که عمق و عرض را ترکیب می‌کنند، مانند شبکه‌های "عمیق و گسترده" تا از مزایای هر دو رویکرد استفاده کنند.

در عمل، انتخاب بین شبکه‌های عمیق تر و گسترده تر به مشکل خاص و منابع موجود بستگی دارد. در اینجا چند دستورالعمل وجود دارد:

1. توابع پیچیده: اگر تابع برای تقریب بسیار پیچیده و سلسله مراتبی باشد، شبکه‌های عمیق تر اغلب کارآمدتر هستند. آنها می‌توانند به طور خودکار روابط پیچیده را یاد بگیرند و به پارامترهای کمتری نیاز دارند.

2. داده‌های محدود: اگر مقدار داده محدودی دارید، شبکه‌های گسترده تر می‌توانند از نظر داده کارآمدتر باشند. آنها می‌توانند انواع مختلفی از ویژگی‌های ساده را به تصویر بکشند و از نصب بیش از حد جلوگیری کنند.

3. منابع محاسباتی: منابع محاسباتی موجود را در نظر بگیرید. آموزش شبکه‌های عمیق تر از نظر محاسباتی پرهزینه است، در حالی که شبکه‌های گسترده تر می‌توانند حافظه فشرده‌ای داشته باشند.

4. منظم سازی: از تکنیک‌های منظم سازی مناسب مانند ترک تحصیل، کاهش وزن و توقف زودهنگام برای جلوگیری از تطبیق بیش از حد در شبکه‌های وسیع تر استفاده کنید.

5. آزمایش تجربی: با معماری های مختلف، از جمله تغییرات در عمق و عرض، آزمایش کنید و از تکنیک هایی مانند اعتبار سنجی متقاطع استفاده کنید تا مشخص کنید کدام معماری برای مشکل خاص شما بهتر عمل می کند.

(د) مزایا:

- نمایش بهبود یافته: مدل های عمیق تر می توانند ویژگی های پیچیده و سلسله مراتبی را در داده ها ثبت کنند. آنها می توانند یاد بگیرند که الگوها و روابط پیچیده را نشان دهند، که برای کارهایی مانند تشخیص تصویر و گفتار، پردازش زبان طبیعی و موارد دیگر بسیار مهم است.
- تعمیم بهتر: شبکه های عمیق اغلب در تعمیم به داده های جدید و دیده نشده بهتر هستند. آنها می توانند ویژگی های انتزاعی و قوی تری را استخراج کنند، که بیش از حد برازش را کاهش داده و توانایی مدل را برای پیش بینی دقیق در طیف وسیعی از ورودی ها افزایش می دهد.
- بیان بالاتر: معماری های عمیق درجه بالاتری از بیان را ارائه می دهند و به آنها اجازه می دهد تا دسته وسیع تری از توابع را تقریب بزنند. این

بدان معنی است که آنها می توانند با طیف وسیع تری از وظایف مقابله کنند و روابط پیچیده تری را مدل کنند.

- یادگیری ویژگی کارآمد: شبکه‌های عمیق می‌توانند به طور خودکار ویژگی‌ها را از داده‌های خام بیاموزند و نیاز به مهندسی ویژگی‌های دستی گسترده را از بین ببرند. این امر به ویژه هنگام برخورد با داده‌های با ابعاد بالا یا بدون ساختار مفید است.

معایب:

- پیچیدگی آموزش: آموزش مدل‌های عمیق پیچیده تر است. شبکه‌های عمیق‌تر ممکن است از مسائلی مانند ناپدید شدن گرادیان‌ها (شیب‌ها بسیار کوچک می‌شوند) یا شیب‌های انفجاری (شیب‌ها بسیار بزرگ می‌شوند) رنج ببرند، که می‌تواند آموزش را چالش‌برانگیز کند.
- منابع محاسباتی: آموزش شبکه‌های عمیق اغلب به منابع محاسباتی قابل توجهی از جمله GPU یا TPU های قدرتمند نیاز دارد. یادگیری عمیق می‌تواند از نظر محاسباتی گران باشد، که ممکن است برای همه پروژه‌ها امکان پذیر نباشد.

- الزامات داده: مدل‌های عمیق معمولاً به داده‌های بیشتری برای تعمیم مؤثر نیاز دارند. اگر داده‌های محدودی دارید، شبکه‌های عمیق‌تر ممکن است مستعد بیش از حد برازش شوند. مدل‌های کم عمق ممکن است در چنین مواردی عملکرد بهتری داشته باشند.

- تنظیم Hyperparameter: شبکه‌های عمیق‌تر دارای هایپرپارامترهای بیشتری برای تنظیم هستند، مانند نرخ یادگیری، اندازه دسته و روش‌های اولیه. یافتن هایپرپارامترهای مناسب می‌تواند وقت گیرتر باشد.

- افزایش خطر بیش‌برازش: در حالی که شبکه‌های عمیق می‌توانند تعمیم‌سازی را بهبود بخشند، مدل‌های بسیار عمیق نیز در صورت عدم تنظیم مناسب در معرض خطر بیش از حد برازش هستند. تکنیک‌هایی مانند ترک تحصیل، عادی‌سازی دسته‌ای و توقف زودهنگام اغلب مورد نیاز است.

به طور خلاصه، افزودن لایه‌های بیشتر به یک شبکه عصبی می‌تواند به بهبود نمایش و تعمیم ویژگی منجر شود و یادگیری عمیق را برای کارهای پیچیده مناسب کند. با این حال، چالش‌های مربوط به پیچیدگی آموزش، منابع

محاسباتی، نیازهای داده، و نیاز به تنظیم و منظم‌سازی موثر فراپارامتر را نیز معرفی می‌کند.

پایان