

## سوال ۱-

- آ. • **دندريت:** دندريت‌ها در يك سلول عصبی، وظیفه انتشار محرک الكتروشیمیایی دریافت شده از دیگر نورون‌ها را به جسم سلولی (سوما) دارند. در واقع دندريت‌ها، معادل بردار ورودی (حاوی مقادیر ضرب شده در وزن‌ها) يك نورون در يك پرسپترون است.
- **سوما:** سوما یا جسم سلولی، به نوعی هسته اصلی يك نورون است و بخش پردازشی خاصی ندارد. از آنجا که ورودی‌های دندريت‌ها به این بخش وارد می‌شوند، می‌توان به نوعی آن را معادل جمع‌کننده (عملگر سیگما) مقادیر ورودی (پس از ضرب در وزن‌ها) در يك نورون يك پرسپترون در نظر گرفت.
- **آکسون:** آکسون در سلول عصبی وظیفه هدایت پالس الكتريکی را دارد و اطلاعات را منتقل می‌کند. در واقع می‌توان آکسون را شبیه به تابع فعال‌سازی و خروجی نورون در يك پرسپترون دانست.
- **سیناپس:** سیناپس، که بخشی از يك نورون نیست، ساختاری است که اجازه انتقال يك سیگنال الكتريکی یا شیمیایی را به نورون بعدی می‌دهد. در واقع در سیناپس، مشابه عمل وزن‌ها در يك شبکه پرسپترون، سیگنال می‌تواند تضعیف یا تقویت شود.
- ب. • **یادگیری با نظارت<sup>۱</sup>:** در این نوع یادگیری، هدف یادگیری تابعی است که ورودی را به خروجی نگاشت می‌کند. این یادگیری به وسیله زوج‌های ورودی و خروجی نمونه و دارای برچسب<sup>۲</sup> انجام می‌شود. از آنجا که شبکه‌های عصبی برای بازنمایی توابع غیرخطی و پیچیده، بسیار مناسب هستند، نقش بسزایی در یادگیری با نظارت دارند.
- **یادگیری بدون نظارت<sup>۳</sup>:** در یادگیری بدون نظارت، داده دارای برچسب وجود ندارد. در نتیجه، الگوریتم یادگیری الگوهای موجود در بین داده‌ها را پیدا می‌کند. نمونه‌های این نوع یادگیری خوشه‌بندی و کاهش بعد هستند. شبکه‌های auto-encoder و استفاده از مقادیر فعال‌سازی لایه‌های مخفی آن‌ها برای خوشه‌بندی و کاهش بعد (استخراج ویژگی‌ها)، در این زمینه پرکاربرد هستند. همچنین معماری‌های شبکه عصبی مانند SOM نیز هستند که می‌توانند به طور مستقیم برای خوشه‌بندی استفاده شوند.

<sup>۱</sup>Supervised Learning<sup>۲</sup>Label<sup>۳</sup>Unsupervised Learning

• **یادگیری تقویتی<sup>۴</sup>**: در یادگیری عمیق، هدف انجام اعمالی توسط یک عامل هوشمند در یک محیط است تا پاداش<sup>۵</sup> دریافتی را حداکثر کند. در این مسائل، با بزرگ شدن محیط و اعمال (افزایش تعداد حالات<sup>۶</sup>)، امکان ذخیره‌سازی بهترین عمل در هر حالت به شکل‌های ساده مانند جدول وجود ندارد. در نتیجه می‌توان از شبکه‌های عصبی (مانند DQN) برای بازنمایی توابع غیرخطی و پیچیده استفاده کرد.

## سوال ۲-

آ. در صورتی که از تابع فعالیت غیرخطی استفاده نکنیم، مستقل از تعداد لایه‌های درونی، می‌توان خروجی را برحسب ترکیب خطی ورودی نوشت. در نتیجه افزودن لایه‌های میانی تنها بار محاسباتی ما را افزایش می‌دهد و بر روی کارایی تاثیری ندارد. به همین دلیل، از توابع فعال‌سازی غیرخطی در هر لایه استفاده می‌شود تا بتوان با افزودن تعداد لایه‌ها توابع پیچیده‌تری را بازنمایی کرد.

ب. افزودن بایاس به نوعی امکان شیفت دادن تابع فعالیت را فراهم می‌کند. این کار به فیت شدن بهتر بر روی داده‌ها کمک می‌کند.

پ. به طور کلی افزایش تعداد نوروها در هر لایه و افزایش تعداد لایه باعث پیچیده‌تر شدن شبکه عصبی می‌شود. یکی اثرات واضح این کار، افزایش بار محاسباتی است. در طرف دیگر، پیچیدگی بیشتر شبکه عصبی، امکان بازنمایی توابع پیچیده‌تر توسط آن و کاهش بایاس (که منجر به underfitting می‌شود) را می‌دهد. اما باید توجه کرد، افزایش عمق می‌تواند روند یادگیری را (به خصوص در لایه‌های ابتدایی) بسیار کند کند. در واقع افزایش تعداد لایه‌ها می‌تواند باعث نزدیک به صفر شدن گرادیان<sup>۷</sup> یا بسیار بزرگ شدن آن<sup>۸</sup> شود که هر دوی آن‌ها یادگیری را مختل می‌کند.

ت. با استفاده از ضریب یادگیری می‌توان اندازه گام را تعیین نمود. معمولاً در ابتدای یادگیری می‌توان از اندازه قدم‌های بزرگ‌تری استفاده کرد و با نزدیک شدن به نقطه کمینه، بر اساس یک قاعده مشخص نرخ یادگیری (و در نتیجه اندازه گام) را کاهش داد. به عنوان مثال، در صورتی که با برداشتن چند قدم، مشاهده کنیم که مقدار تابع هزینه کاهش پیدا نکرده است، می‌توان نرخ

<sup>4</sup>Reinforcement Learning

<sup>5</sup>Reward

<sup>6</sup>States

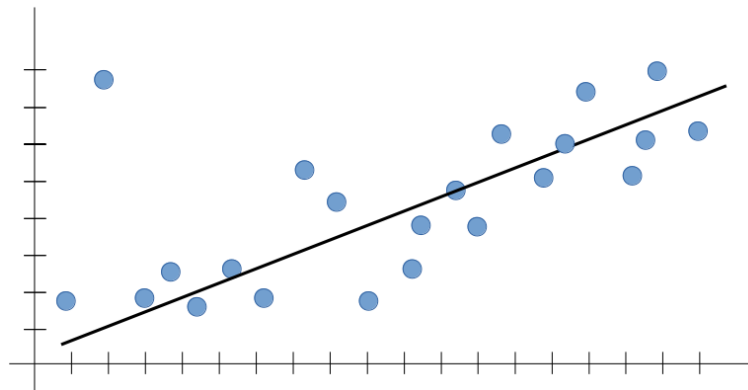
<sup>7</sup>Vanishing gradient

<sup>8</sup>Exploding gradient

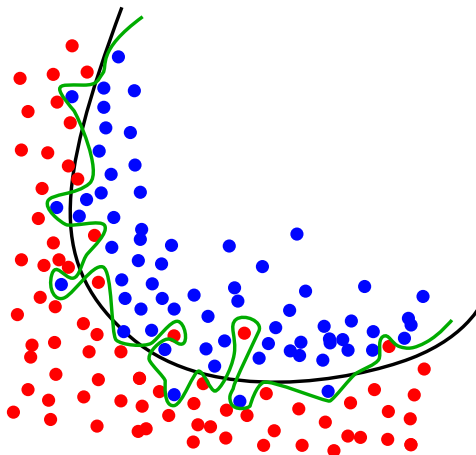
یادگیری را کاهش داد (مثلا در ۰/۹ ضرب کنیم).

### سوال ۳-

آ. آندرفیتینگ زمانی رخ می دهد که یک مدل نتواند به میزان کافی بر روی ساختار داده ها یادگیری داشته باشد. به عنوان مثال در شکل زیر، مدل خطی نمی تواند پیش بینی مناسبی بر روی داده ها داشته باشد.

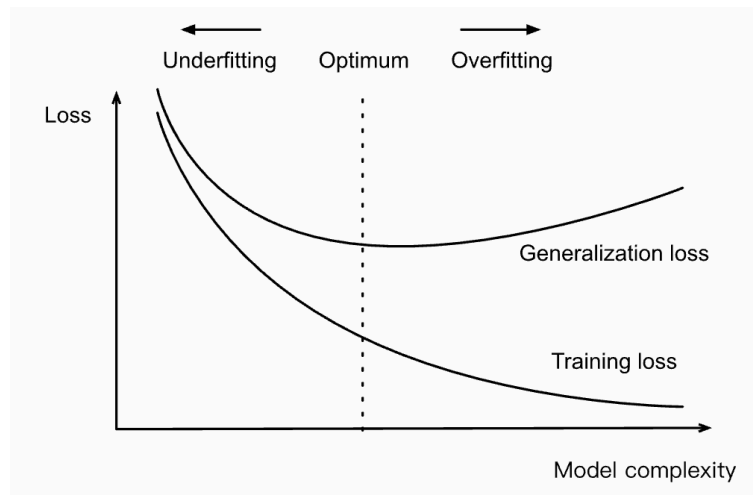


در طرف مقابل، اورفیتینگ زمانی رخ می دهد که مدل بسیار به داده های خاصی (داده های آموزش) وابسته شود و در مواجهه با نمونه های جدید نتواند به خوبی آن ها را پیش بینی کند. در واقع مدل عمومیت نخواهد داشت. در شکل زیر، خط سبز رنگ نشان دهنده یک مدل اورفیت شده است که عمومیت مناسبی ندارد. در حالی که خط سیاه رنگ چنین نیست و اگر چه خطای بیشتری روی داده های آموزشی دارد، اما در مواجهه با نمونه های جدید دقت بالاتری دارد.



هر دوی این اتفاقات، دقت پیش بینی را کاهش می دهد.

ب. برای این کار ابتدا داده‌ها را به دو مجموعه train و test تقسیم می‌کنیم. همانطور که در شکل زیر می‌بینیم، در صورتی که میزان loss هر دو مجموعه train و test عدد بالایی بود، دچار آندرفیتینگ هستیم. در طرف مقابل، در صورتی که loss مجموعه train کاهش پیدا کند، اما مجموعه test افزایش پیدا کند (بین این دو فاصله بیفتد) دچار اورفیتینگ هستیم و مدل عمومیت ندارد.



زمانی که هر دو نمودار loss کاهش داشته باشند (شکل زیر)، یعنی این دو مشکل وجود ندارد.



پ. برای رفع مشکل آندرفیتینگ، یکی از راه‌ها آموزش دادن مدل به مدل طولانی‌تر است. اما باید توجه کرد این کار لزوماً نمی‌تواند مشکل را رفع کند؛ به عنوان مثال مدل خطی در **مورد آ**. امکان بازنمایی با دقت بالای داده‌های مسئله را ندارد. در نتیجه می‌توان از مدل‌های پیچیده‌تر (مثلاً مدل چندجمله‌ای به جای خطی یا استفاده از شبکه عصبی عمیق‌تر) استفاده کرد.

ت. • **روش regularization**: در این روش در کنار تابع loss مورد استفاده (مثلاً cross-entropy)، قدر مطلق (L1 regularization) یا مجذور (L2 regularization) وزن‌ها را با آن جمع می‌کنیم؛ یعنی:

$$\mathcal{L}_{\mathcal{R}} = \mathcal{L} + \frac{\lambda}{2} \sum_{l=1}^L ||W^{[l]}||^2$$

که در آن  $\lambda$  پارامتر regularization است و نسبت اهمیت بخش cross-entropy و regularization را مشخص می‌کند. با این کار در واقع شبکه را برای بزرگ شدن وزن‌های آن جریمه می‌کنیم و در نتیجه وزن‌های شبکه مقادیر کوچک‌تری به خود می‌گیرند. از آنجا که زمانی که اورفیتینگ رخ می‌دهد، برخی از وزن‌های شبکه عصبی مقادیر بزرگی به خود می‌گیرند و تاثیر بالایی بر روی خروجی می‌گذارند، انجام این کار از اورفیتینگ جلوگیری می‌کند.

• **روش dropout:** در این روش، با انتخاب یک پارامتر بین صفر و یک، در آموزش شبکه عصبی، برای هر نورون در هر اجرای backprop، یک عدد تصادفی بین صفر و یک انتخاب می‌شود و بسته به کوچک‌تر یا بزرگ‌تر بودن آن از پارامتر انتخابی، تاثیر یا عدم تاثیرگذاری آن نورون در آن مرحله مشخص می‌شود. این کار به نوعی باعث داشتن یک شبکه عصبی ساده‌تر در هر مرحله و همچنین پخش شدن مقادیر وزن‌ها روی نورون‌های مختلف می‌شود. در نتیجه از اورفیتینگ جلوگیری می‌شود.

#### سوال ۴-

آ. طبق قاعده مشتق زنجیره‌ای داریم (مقدار هر یک از نورون‌ها پیش از اعمال تابع فعال‌سازی با حرف کوچک نشان داده شده است):

$$\frac{\partial \mathcal{L}}{\partial u_0} = \frac{\partial \mathcal{L}}{\partial O_1} \cdot \frac{\partial O_1}{\partial o_1} \cdot \frac{\partial o_1}{\partial B_0} \cdot \frac{\partial B_0}{\partial b_0} \cdot \frac{\partial b_0}{\partial A_1} \cdot \frac{\partial A_1}{\partial a_1} \cdot \frac{\partial a_1}{\partial u_0}$$

مشتق تابع  $g(x) = \sigma(x)$  برابر است با:

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}} = \sigma(x) \cdot (1-\sigma(x))$$

حال هر یک از مشتق‌های جزئی را به طور جداگانه محاسبه می‌کنیم.

$$\frac{\partial \mathcal{L}}{\partial O_1} = 2(O_1 - y_t)$$

$$\frac{\partial O_1}{\partial o_1} = \sigma(o_1) (1 - \sigma(o_1))$$

$$\frac{\partial o_1}{\partial B_0} = t_0$$

$$\frac{\partial B_0}{\partial b_0} = \sigma(b_0) (1 - \sigma(b_0))$$

$$\frac{\partial b_0}{\partial A_1} = v_1$$

$$\frac{\partial A_1}{\partial a_1} = \sigma(a_1) (1 - \sigma(a_1))$$

$$\frac{\partial a_1}{\partial u_0} = X_0$$

پس در نهایت با جایگذاری مشتق‌ها، خواهیم داشت:

$$\frac{\partial \mathcal{L}}{\partial u_0} = 2(O_1 - y_t) \cdot \sigma(o_1) (1 - \sigma(o_1)) \cdot t_0 \cdot \sigma(b_0) (1 - \sigma(b_0)) \cdot v_1 \cdot \sigma(a_1) (1 - \sigma(a_1)) \cdot X_0$$

ب. با انجام دادن محاسبات به شکل برداری، به ترتیب هر لایه را حساب می‌کنیم (مقدار هر یک از نورون‌ها پیش از اعمال تابع فعال‌سازی با حرف کوچک نشان داده شده است).

$$a = \begin{bmatrix} a_0 \\ a_1 \end{bmatrix} = \begin{bmatrix} b_0 & w_0 & w_1 & w_2 \\ b_1 & u_0 & u_1 & u_2 \end{bmatrix} \begin{bmatrix} 1 \\ X_0 \\ X_1 \\ X_2 \end{bmatrix} = \begin{bmatrix} -0.3 & 0.1 & 0.5 & 0.7 \\ 0.5 & 0.6 & 0.9 & 0.1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 1.2 \end{bmatrix}$$

$$A = \sigma(a) = \sigma \left( \begin{bmatrix} 0.5 \\ 1.2 \end{bmatrix} \right) = \begin{bmatrix} 0.62 \\ 0.76 \end{bmatrix}$$

$$b = \begin{bmatrix} b_0 \end{bmatrix} = \begin{bmatrix} b_2 & v_0 & v_1 \end{bmatrix} \begin{bmatrix} 1 \\ A_0 \\ A_1 \end{bmatrix} = \begin{bmatrix} -0.4 & 0.2 & 0.6 \end{bmatrix} \begin{bmatrix} 1 \\ 0.62 \\ 0.76 \end{bmatrix} = \begin{bmatrix} 0.18 \end{bmatrix}$$

$$B = \sigma(b) = \sigma \left( \begin{bmatrix} 0.18 \end{bmatrix} \right) = \begin{bmatrix} 0.54 \end{bmatrix}$$

$$o = \begin{bmatrix} o_0 \end{bmatrix} = \begin{bmatrix} b_3 & t_0 \end{bmatrix} \begin{bmatrix} 1 \\ B_0 \end{bmatrix} = \begin{bmatrix} 0.2 & 0.8 \end{bmatrix} \begin{bmatrix} 1 \\ 0.54 \end{bmatrix} = \begin{bmatrix} 0.63 \end{bmatrix}$$

$$O_1 = \sigma(0.63) = 0.65$$

میزان خطا (تابع cost) برابر است با:

$$\text{cost} = (0.65 - 0)^2 = 0.42$$

## سوال ۵-

آ. بر خلاف FCNN که تمام نورون‌ها در یک لایه به تمام نورون‌ها در لایه بعدی متصل هستند، در شبکه‌های پیچشی، تنها بخش کوچک و نزدیکی از هر لایه، به لایه بعد متصل می‌شود (اتصال محلی).

همچنین از آنجا که هر فیلتر (پارامترها) بر روی کل نورون‌های لایه قبل جابجا می‌شود، این فیلتر ویژگی‌ها را از تمام قسمت‌های لایه قبل استخراج می‌کند (به اشتراک‌گذاری پارامترها). هر دوی این ویژگی‌ها، یعنی اتصال بخش‌های کوچکی از لایه‌های مجاور و به اشتراک‌گذاری ویژگی‌های یاد گرفته شده (مثلا درک خطوط، اشکال و ...)، باعث کارایی مناسب CNN برای داده‌های عکسی می‌شوند.

ب. طبق فرمول زیر، طول و عرض خروجی برابر است با:

$$n' = \left\lfloor \frac{n + 2p - f}{s} \right\rfloor + 1 = \left\lfloor \frac{64 + 2 \times 2 - 5}{3} \right\rfloor + 1 = 22$$

از آنجا که ۶ فیلتر داریم، پس ابعاد خروجی برابر  $22 \times 22 \times 6$  است.

پ. هر فیلتر با احتساب بایاس،  $5 \times 5 \times 3 + 1 = 76$  پارامتر دارد. پس در مجموع  $76 \times 6 = 456$  پارامتر داریم.

## سوال ۶-

آ. در صورتی که مدل، کلاس داده ورودی را به درستی و با اختلاف بیش از یک نسبت به بقیه کلاس‌ها پیش‌بینی کند، این تابع مقدار صفر را برمی‌گرداند. از آنجا که در این تابع ورودی‌ها با صفر max شده‌اند، پس کمترین مقدار این تابع صفر خواهد بود. بیشینه مقدار این تابع اما مشخص نیست و هر مقدار نامنفی را می‌تواند اختیار کند.

ب. در این صورت عبارت  $s_j - s_{yi} + 1$  برای تمام نوروهای خروجی تقریباً برابر  $0 - 0 + 1 = 1$  می‌شود و خروجی max نیز تقریباً برابر یک می‌شود. در نتیجه با توجه به اینکه  $N - 1$  (به دلیل جمع نکردن کلاس ground-truth) جمع انجام می‌شود، مقدار  $\mathcal{L}_i$  تقریباً برابر  $N - 1$  خواهد شد.

پ. در این صورت یک واحد به تمامی  $\mathcal{L}_i$ ‌ها و همچنین تابع cost اضافه می‌شود. در نتیجه مقدار کمینه تابع به جای صفر، برابر یک خواهد بود و تاثیری بر روی بهینه‌سازی این تابع و یادگیری پارامترها نخواهد داشت.