



دانشگاه صنعتی شریف

دانشکده مهندسی کامپیوتر

طراحان: سپهرغیاث، امیرکوشان فتاح، شایان بقایی نژاد، آرش ضیائی، رادین چراغی، علیرضا ملک حسینی

## هوش مصنوعی

بهار ۱۴۰۴

استاد: احسان تن قطاری

تمرین چهارم

یادگیری ماشین و شبکه عصبی

مهلت ارسال: ۴ خرداد

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- در طول ترم امکان ارسال با تاخیر پاسخ همه‌ی تمارین سقف ۴ روز و در مجموع ۱۰ روز، وجود دارد. پس از گذشت این مدت، پاسخ‌های ارسال شده پذیرفته نخواهند بود. همچنین، به ازای هر ساعت تأخیر غیر مجاز نیم درصد از نمره‌ی تمرین کم خواهد شد.
- هم‌کاری و هم‌فکری شما در انجام تمرین مانعی ندارد اما پاسخ‌های ارسال‌شده هر کس حتماً باید توسط خود او نوشته شده باشد.
- در صورت هم‌فکری و یا استفاده از هر منابع خارج درسی، نام هم‌فکران و آدرس منابع مورد استفاده برای حل سوال مورد نظر را ذکر کنید.
- لطفاً تصویری واضح از پاسخ‌های سوالات نظری بارگذاری کنید. در غیر این صورت پاسخ شما تصحیح نخواهد شد.

### سوالات (۱۰۰ نمره)

۱. (۵ نمره) درستی یا نادرستی عبارت‌های زیر را با ذکر دلیل مشخص کنید.
  - (آ) افزایش پیچیدگی یک مدل رگرسیون، همواره منجر به کاهش خطای مدل در داده‌های آموزش و افزایش خطای مدل در برآورد داده تست می‌شود.
  - (ب) اگر bias مدل ما زیاد است، اضافه کردن تعداد داده‌های آموزش کمک زیادی به کم کردن bias نمی‌کند.
  - (ج) از آنجا که طبقه بندی حالت خاصی از رگرسیون است، رگرسیون لجیستیک نیز حالت خاصی از رگرسیون خطی می‌باشد.
  - (د) در صورتی که یک نمونه از داده‌های آموزشی توسط الگوریتم پرسپترون اشتباه کلاس بندی شود، پس از بروزرسانی وزن‌های مدل با استفاده از همان نمونه، الگوریتم آنرا به درستی کلاس بندی می‌کند.

حل:

- (آ) غلط، چرا که اگر مدل در حالت underfitting باشد (یعنی بیش از حد ساده باشد و نتواند الگوی موجود در داده را یاد بگیرد)، افزایش پیچیدگی می‌تواند هم خطای آموزش و هم خطای تست را کاهش دهد، چون مدل بهتر با الگوی واقعی داده منطبق می‌شود. تنها زمانی که مدل به مرحله overfitting برسد، پیچیدگی بیشتر باعث افزایش خطای تست می‌شود.
- (ب) غلط، در کل بایاس بالا معمولاً ناشی از ساده بودن مدل است. (underfitting) در این حالت، حتی با اضافه کردن داده‌های بیشتر، مدل همچنان قادر به یادگیری الگوی پیچیده مسئله نیست. کاهش بایاس نیازمند افزایش ظرفیت مدل (مثلاً استفاده از مدل‌های غیرخطی‌تر یا پیچیده‌تر) است و افزایش داده احتمالاً کمک خاصی به آن نمی‌کند. مگر اینکه تعداد داده‌ها قبل از این بیش از حد کم و پراکنده بوده و الگوی موجود در توزیع اصلی داده‌ها در سмпل اولیه یافت نمی‌شود (و به نحوی الگوی پیچیده‌تری ساخته شده). در این حالت ممکن است (به احتمال کمی) بایاس کم شود.

(ج) غلط، هرچند از نظر ساختار تابع، رگرسیون لجیستیک از یک ترکیب خطی از ویژگی‌ها استفاده می‌کند، اما خروجی آن احتمال (بین ۰ و ۱) است که از طریق تابع سیگموئید حاصل می‌شود، نه یک مقدار عددی پیوسته مانند رگرسیون خطی. بنابراین، رگرسیون لجیستیک را نمی‌توان به طور کامل «حالت خاصی از رگرسیون خطی» دانست. بلکه می‌توان گفت از ترکیب خطی در مرحله تصمیم‌گیری بهره می‌برد ولی ماهیت آن متفاوت است.

(د) غلط، به‌روزرسانی وزن‌ها لزوماً باعث تصحیح فوری طبقه‌بندی نمونه نمی‌شود. ممکن است چندین تکرار لازم باشد.

یک مثال نقض می‌زنیم: فرض کنید وزن اولیه و داده به صورت زیر باشد

$$\mathbf{w} = \begin{bmatrix} 5 \\ 5 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad y = -1$$

ابتدا مدل یک پیش‌بینی اشتباه ارائه می‌دهد.

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x}) = \text{sign}(5 \cdot 1 + 5 \cdot 1) = \text{sign}(10) = +1$$

اشتباه طبقه‌بندی شده است.

حال به‌روزرسانی وزن‌ها انجام می‌شود.

$$\mathbf{w}_{\text{new}} = \mathbf{w} + y \cdot \mathbf{x} = \begin{bmatrix} 5 \\ 5 \end{bmatrix} + (-1) \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \end{bmatrix}$$

بعد از به‌روزرسانی پیش‌بینی را دوباره انجام می‌دهیم.

$$f(\mathbf{x}) = \text{sign}(4 \cdot 1 + 4 \cdot 1) = \text{sign}(8) = +1$$

می‌بینیم که همچنان اشتباه طبقه‌بندی شده است.

۲. (۱۲ نمره) شما تصمیم گرفته‌اید معلم شوید. تنها مشکلی که وجود دارد این است که نمی‌خواهید زمان زیادی را صرف تصحیح مقاله‌ها کنید، بنابراین تصمیم می‌گیرید همه آن‌ها را با یک طبقه‌بند خطی نمره‌دهی کنید. طبقه‌بند شما تعداد کلمات ۷ حرفی ( $f_7$ ) و ۸ حرفی ( $f_8$ ) در یک مقاله را در نظر می‌گیرد و سپس براساس این دو عدد به مقاله نمره A یا F می‌دهد. شما چهار مقاله نمره‌دار برای یادگیری در اختیار دارید:

نمره (A (+), F (-))	تعداد کلمات ۸ حرفی ( $f_8$ )	تعداد کلمات ۷ حرفی ( $f_7$ )
A	۲	۳
F	۱	۱
A	۳	۴
F	۲	۰

(آ) ابتدا با ذکر دلیل مشخص کنید که داده‌های آموزشی با ویژگی‌های داده‌شده به‌صورت خطی جداپذیر هستند یا نه.

(ب) شما تصمیم می‌گیرید الگوریتم پرسپترون را اجرا کنید و با خوش‌بینی نسبت به توانایی‌های نگارش مقاله‌ی دانش‌آموزان، بردار وزن خود را به صورت (۰, ۰, ۱) مقداردهی اولیه می‌کنید. اگر امتیاز خروجی طبقه‌بند شما بیشتر از ۰ باشد، نمره‌ی A می‌دهد؛ اگر ۰ یا کمتر باشد، نمره‌ی F می‌دهد. بردار وزن حاصل را پس از مشاهده‌ی اولین و دومین مثال آموزشی با استفاده از الگوریتم پرسپترون بنویسید.

(ج) (۴ نمره) برای هر یک از قوانین تصمیمگیری زیر، مشخص کنید آیا بردار وزنی وجود دارد که آن قانون را نمایش دهد. اگر پاسخ «بله» است، بردار وزنی مربوطه را بنویسید.

- مقاله نمره A میگیرد اگر و تنها اگر  $(f_7 + f_8 \geq 7)$  برقرار باشد.
- مقاله نمره A میگیرد اگر و تنها اگر  $(f_7 \geq 5$  و  $f_8 \geq 4)$  برقرار باشد.
- مقاله نمره A میگیرد اگر و تنها اگر  $(f_7 \geq 5$  یا  $f_8 \geq 4)$  برقرار باشد.
- مقاله نمره A میگیرد اگر و تنها اگر بین ۴ تا ۶ (شامل هر دو) کلمه ۷ حرفی و بین ۳ تا ۵ کلمه ۸ حرفی داشته باشد.

حل.

(آ) درست. این داده‌ها با یک مرز خطی قابل جدا شدن هستند. به عنوان مثال، بردار وزن  $(2, 5, 1)$  می‌تواند همه نقاط را به درستی طبقه‌بندی کند.

(ب) وزن اولیه برابر است با

$$(1, 0, 0)$$

پس از دیدن نمونه اول با بردار ویژگی  $(1, 2, 1)$  و برچسب مثبت،  $(A)$  حاصل ضرب داخلی برابر است با:

$$1 \cdot 1 + 0 \cdot 2 + 0 \cdot 1 = 1 > 0$$

پس پیش‌بینی درست بوده و وزن تغییر نمی‌کند:

$$(1, 0, 0)$$

پس از دیدن نمونه دوم با بردار ویژگی  $(1, 0, 2)$  و برچسب منفی،  $(F)$  حاصل ضرب داخلی برابر است با:

$$1 \cdot 1 + 0 \cdot 0 + 0 \cdot 2 = 1 > 0$$

پیش‌بینی اشتباه بوده، پس به‌روزرسانی انجام می‌دهیم:

$$(1, 0, 0) + (-1) \cdot (1, 0, 2) = (0, 0, -2)$$

(ج) • بله. شرط  $f_7 + f_8 \geq 7$  یک مرز خطی ایجاد می‌کند و می‌توان آن را با بردار وزن  $(1, 1, -7)$  نمایش داد.

• خیر. شرط  $f_8 \geq 4$  و  $f_7 \geq 5$  یک منطقه اشتراکی غیرخطی تعریف می‌کند و قابل نمایش با یک بردار وزن نیست.

• خیر. شرط  $f_8 \geq 4$  یا  $f_7 \geq 5$  یک اجتماع از دو ناحیه جداگانه است و قابل نمایش با یک مرز خطی نیست.

• خیر. داشتن تعداد مشخصی از کلمات بین دو بازه، یک ناحیه بسته و محدود (مستطیلی) تعریف می‌کند که نمی‌توان آن را با یک مرز خطی نمایش داد.

۳. (۱۵ نمره) فرض کنید یک مجموعه داده با سه ویژگی ورودی دودویی به صورت زیر  $A, B, C$  و یک ویژگی خروجی دودویی  $Y$  در اختیار داریم.

ویژگی‌های ورودی مقادیر  $\{0, 1\}$  می‌گیرند، در حالی که  $Y$  مقادیر  $\{\text{درست}, \text{نادرست}\}$  را می‌گیرد.

$Y$	$C$	$B$	$A$
درست	۱	۱	۰
درست	۰	۱	۱
نادرست	۱	۰	۱
نادرست	۱	۱	۱
درست	۱	۱	۰
درست	۰	۰	۰
نادرست	۱	۱	۰
نادرست	۱	۰	۱
درست	۰	۱	۰
درست	۱	۱	۱

با توجه به اطلاعات ذکر شده، به موارد زیر پاسخ دهید.

(آ) Naive Bayes رکورد  $(A = ۱, B = ۱, C = ۰)$  را چگونه طبقه‌بندی می‌کند؟

(ب) Naive Bayes رکورد  $(A = ۰, B = ۰, C = ۱)$  را چگونه طبقه‌بندی می‌کند؟

(ج) Naive Bayes رکورد  $(A = ۰, B = ۰, C = ۰)$  را با استفاده از Laplace smoothing(k=1) چگونه طبقه‌بندی می‌کند؟

(د) آیا افزودن تنها یک رکورد می‌تواند باعث تغییر طبقه‌بندی رکورد  $(A = ۱, B = ۰, C = ۱)$  شود؟

حل.

(آ) «درست»  $\rightarrow (1,1,0)$

$$P_T = P(Y=\text{درست}) P(A=۱|T) P(B=۱|T) P(C=۰|T) = \frac{6}{10} \times \frac{2}{6} \times \frac{5}{6} \times \frac{3}{6} = \frac{1}{12},$$

$$P_F = P(Y=\text{نادرست}) P(A=۱|F) P(B=۱|F) P(C=۰|F) = \frac{4}{10} \times \frac{3}{4} \times \frac{2}{4} \times \frac{0}{4} = 0.$$

چون  $P_T > P_F$ ، برچسب «درست» برگردانده می‌شود.

(ب) «نادرست»  $\rightarrow (0,0,1)$

$$P_T = \frac{6}{10} \times \frac{4}{6} \times \frac{1}{6} \times \frac{3}{6} = 0.0333\dots,$$

$$P_F = \frac{4}{10} \times \frac{1}{4} \times \frac{2}{4} \times \frac{4}{4} = 0.05.$$

از آنجا که  $P_F > P_T$  است، «نادرست» پیش‌بینی می‌شود.

(ج) «درست»  $\rightarrow (0,0,0)$  فرمول هموارسازی لاپلاس برای متغیرهای دودویی:

$$P(X=x|Y=y) = \frac{\text{count}(X=x, Y=y) + 1}{N_y + 2}, \quad P(Y=y) = \frac{\text{count}(Y=y) + 1}{N + 2}.$$

با  $N = ۱۰$  و  $|Y| = ۲$ :

$$P(Y=T) = \frac{6+1}{10+2} = \frac{7}{12},$$

$$P(Y=F) = \frac{4+1}{10+2} = \frac{5}{12},$$

$$P(A=\bullet|T) = \frac{4+1}{6+2} = \frac{5}{8},$$

$$P(B=\bullet|T) = \frac{1+1}{6+2} = \frac{2}{8},$$

$$P(C=\bullet|T) = \frac{3+1}{6+2} = \frac{4}{8},$$

$$P(A=\bullet|F) = \frac{1+1}{4+2} = \frac{2}{6},$$

$$P(B=\bullet|F) = \frac{2+1}{4+2} = \frac{3}{6},$$

$$P(C=\bullet|F) = \frac{0+1}{4+2} = \frac{1}{6}.$$

$$P_T = \frac{7}{12} \times \frac{5}{8} \times \frac{2}{8} \times \frac{4}{8} = \frac{35}{768} \approx 0.0456,$$

$$P_F = \frac{5}{12} \times \frac{2}{6} \times \frac{3}{6} \times \frac{1}{6} = \frac{5}{432} \approx 0.0116.$$

از آنجا که  $P_T > P_F$ ، برچسب نهایی «درست» است. (هموارسازی، احتمال صفر طرف «نادرست» را برطرف می‌کند اما نتیجه تغییری نمی‌کند.)

(د) «خیر»  $(1,0,1) \rightarrow$  در وضعیت فعلی:

$$P_T = \frac{6}{10} \times \frac{2}{6} \times \frac{1}{6} \times \frac{3}{6} = 0.0167,$$

$$P_F = \frac{4}{10} \times \frac{3}{4} \times \frac{2}{4} \times \frac{4}{4} = 0.15.$$

بعد از افزودن داده:

$$P_T = \frac{7}{11} \times \frac{3}{7} \times \frac{2}{7} \times \frac{4}{7} = 0.0445,$$

$$P_F = \frac{5}{11} \times \frac{3}{5} \times \frac{2}{5} \times \frac{4}{5} = 0.0873.$$

افزودن تنها یک نمونه (با هر برچسب یا ویژگی) حداکثر می‌تواند  $P_T$  را تا حدود  $0.445$   $\approx$  افزایش داده یا  $P_F$  را تا حدود  $0.873$   $\approx$  کاهش دهد؛  
 با این حال همچنان  $P_F > P_T$  باقی می‌ماند، در نتیجه پاسخ نهایی «خیر» است.

۴. (۱۶ نمره) به پرسش‌های زیر با استفاده از مجموعه داده سوال قبل و جدول آنتروپی زیر پاسخ دهید.

Specific Conditional Entropies		
$H(Y A = 1, C = 0) = 0.00$	$H(Y A = 0, B = 0) = 0.00$	$H(Y A = 0) = 0.72$
$H(Y A = 1, C = 1) = 0.81$	$H(Y A = 0, B = 1) = 0.81$	$H(Y A = 1) = 0.97$
$H(Y B = 0, C = 0) = 0.00$	$H(Y A = 1, B = 0) = 0.00$	$H(Y B = 0) = 0.92$
$H(Y B = 0, C = 1) = 0.00$	$H(Y A = 1, B = 1) = 0.92$	$H(Y B = 1) = 0.86$
$H(Y B = 1, C = 0) = 0.00$	$H(Y A = 0, C = 0) = 0.00$	$H(Y C = 0) = 0.00$
$H(Y B = 1, C = 1) = 0.97$	$H(Y A = 0, C = 1) = 0.92$	$H(Y C = 1) = 0.99$

(آ) کدام یک از A, B, C بیشترین information gain را دارد.

(ب) درخت تصمیم را بدون استفاده از عملیات هرس برای این سوال رسم کنید. در هنگام ساخت درخت، اگر در مرحله‌ای به دو ویژگی یکسان برای تقسیم برخوردید، ویژگی را براساس ترتیب حروف الفبا انتخاب کنید. همچنین برای برچسب‌گذاری برگ‌ها، از True در صورت تساوی استفاده کنید.

(ج) خروجی درخت تصمیم شما برای ورودی‌های زیر چیست؟

$$(A = 0, B = 0, C = 1) \quad (1)$$

$$(A = 1, B = 0, C = 0) \quad (2)$$

(د) اگر تمامی رتوس به جز ریشه را از درخت تصمیم هرس کنید و فرض کنید در صورت تساوی True را انتخاب کنید. خروجی برای ورودی‌های زیر چه خواهد بود؟

$$(A = 0, B = 0, C = 1) \quad (1)$$

$$(A = 1, B = 0, C = 0) \quad (2)$$

حل.

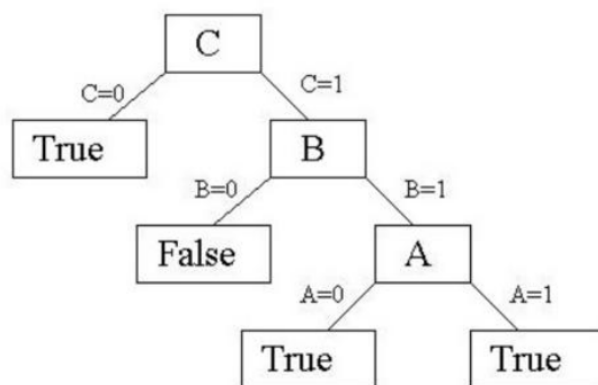
(آ)

$$H(Y|A) = P(A = 0) * H(Y|A = 0) + P(A = 1) * H(Y|A = 1) = \frac{5}{10} * 0.72 + \frac{5}{10} * 0.97 = 0.845$$

$$H(Y|B) = P(B = 0) * H(Y|B = 0) + P(B = 1) * H(Y|B = 1) = \frac{3}{10} * 0.92 + \frac{7}{10} * 0.86 = 0.878$$

$$H(Y|C) = P(C = 0) * H(Y|C = 0) + P(C = 1) * H(Y|C = 1) = \frac{3}{10} * 0.00 + \frac{7}{10} * 0.99 = 0.693$$

information gain زمانی بیشینه می‌شود که  $X = C$



(ب) توجه کنید که در حالتی که  $C = 1$ ، ویژگی اولی که انتخاب می‌شود  $B$  است نه  $A$ ، چرا که:

$$H(Y|A, C=1) = P(A=0|C=1)H(Y|A=0, C=1) + P(A=1|C=1)H(Y|A=1, C=1)$$

$$= \frac{3}{7} * 0.92 + \frac{4}{7} * 0.81 = 0.857$$

که بیشتر از:

$$H(Y|B, C=1) = P(B=0|C=1)H(Y|B=0, C=1) + P(B=1|C=1)H(Y|B=1, C=1)$$

$$= \frac{2}{7} * 0.00 + \frac{5}{7} * 0.97 = 0.693$$

(ج) False •

True •

(د) False •

True •

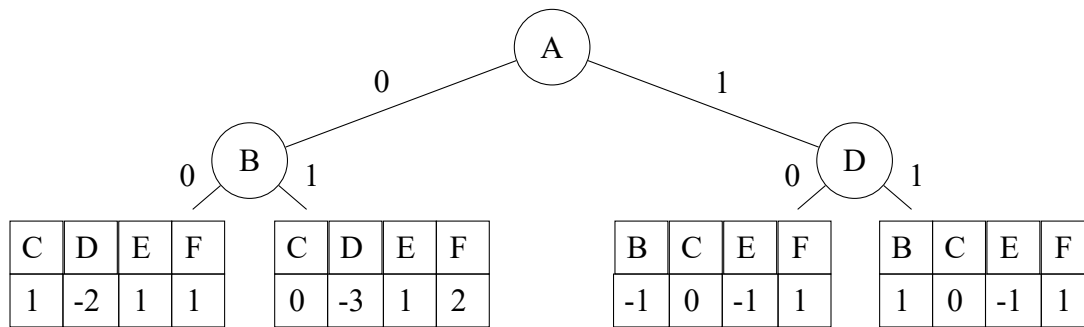
۵. (۱۲ نمره)

نیکا که به دنبال دسته‌بندی کننده‌ی مناسب خودش می‌گشت، الگوریتم جدیدی به نام «درخت پرسپترون» به ذهنش رسید که ویژگی‌هایی را از هر دو مدل درخت تصمیم و الگوریتم پرسپترون را باهم آمیخته می‌کند. درخت پرسپترون مشابه درخت تصمیم است، با این تفاوت که در هر برگ درخت به جای رای اکثریت<sup>۱</sup> از یک مدل پرسپترون استفاده می‌شود.

برای ساخت درخت پرسپترون، نخستین گام دنبال کردن روند عادی یک الگوریتم درخت تصمیم و جدا کردن بر اساس ویژگی‌هاست تا زمانی که به حداکثر عمق مجاز درخت برسیم. پس از رسیدن به این عمق، در هر برگ یک مدل پرسپترون بر پایه‌ی ویژگی‌های استفاده نشده در آن شاخه یاد گرفته می‌شود. در نتیجه با داشتن یک نمونه‌ی جدید، ابتدا نمونه را وارد درخت تصمیم خود می‌کنیم و با پیش‌روی به یکی از برگ‌ها می‌رسیم. سپس مدل پرسپترون آن برگ را بر روی ویژگی‌های بررسی نشده انجام می‌دهیم و دسته‌بندی می‌کنیم.

<sup>۱</sup> Majority Voting

فرض کنید که یک dataset با ۶ ویژگی دودویی<sup>۲</sup>  $\{A, B, C, D, E, F\}$  و دو دسته‌ی خروجی  $\{1, -1\}$  دارید. یک درخت پرسپترون با عمق ۲ برای این dataset در شکل زیر داده شده است. وزن‌های پرسپترون در هر برگ داده شده است. فرض کنید مقدار bias برای هر برگ  $b = 1$  است.



(آ) خروجی درخت پرسپترون برای نمونه‌ی  $x = \{1, 1, 0, 1, 0, 1\}$  چیست؟

(ب) درستی یا نادرستی عبارت‌های زیر را درباره‌ی درخت پرسپترون مشخص کنید.

۱. مرز تصمیم‌گیری درخت پرسپترون همواره خطی است.
۲. برای مقادیر کوچک حداکثر عمق (۲-۳)، احتمال underfit شدن درخت پرسپترون بیشتر از درخت تصمیم است.

(ج) فرض کنید  $D$  که یک dataset است، به ما داده شده است. از دو ساختار درخت تصمیم متفاوت برای یادگیری این dataset استفاده می‌کنیم. یک بار از Information Gain و بار دیگر از Training Error Rate برای جدا کردن نمونه‌ها استفاده می‌کنیم (در حالت دوم در هر گام بر اساس ویژگی‌ای تقسیم‌بندی می‌کنیم که کم‌ترین Training Error را می‌دهد). با فرض این که هر دو درخت تا رسیدن به Training Error صفر آموزش داده می‌شوند، کدام موارد زیر درست هستند؟

- هر دو درخت در ریشه بر پایه‌ی ویژگی یکسانی تقسیم‌بندی می‌کنند.
- هر دو درخت عمق یکسان خواهند داشت.
- هر دو درخت برای هر Data Point در  $D$  خروجی یکسان خواهند داشت.
- هر دو درخت برای هر ورودی، خروجی یکسان خواهند داشت.
- هر دو درخت از تمام ویژگی‌های درون Dataset استفاده خواهند کرد.

حل.

(آ) با توجه به این که  $A = D = 1$  نمونه به راست‌ترین برگ درخت می‌رود که در آن جا

$$\text{out} = 1 \times 1 + 0 \times 0 + (-1) \times 0 + 1 \times 1 + 1 = 3 \Rightarrow \text{sign}(3) = 1$$

پس خروجی برابر ۱ است.

(ب) ۱. نادرست - درخت پرسپترون در ساختار خود درخت تصمیم را نیز دارد که یک Classifier غیر خطی است.

۲. نادرست - به طور کلی درخت پرسپترون ساختار پیچیده‌تری از درخت تصمیم با عمق یکسان دارد و در نتیجه بیشتر تمایل به overfit شدن دارد.

(ج) با توجه به این که می‌دانیم در هر دو  $\text{Error Training} = 0$  نتیجه می‌گیریم هر دو تمام نمونه‌های درون dataset را درست دسته‌بندی می‌کنند. پس مورد سوم درست است. هیچ کدام از موارد دیگر لزوماً برقرار نیستند.



۶. (۲۰ نمره) یکی از روش‌های منظم‌سازی در مسائل رگرسیون خطی، روش Lasso است. در این روش، نرم L1 وزن‌های مدل در تابع خطا وارد می‌شود. این کار باعث می‌شود که پاسخ نهایی مسئله به صورت پراکنده‌تری (sparse) باشد. در این مسئله، خواهیم دید که چگونه جمله‌ی نرم L1 منجر به افزایش پراکندگی می‌شود. هدف پراکنده کردن جواب نهایی (تنک کردن آن) یا همون sparse کردن آن بهینه کردن فضای ذخیره سازی داده ها، ساده کردن محاسبات بعدی و جلوگیری از بیش برازش (overfitting) هست.  $\mathbf{X} \in \mathbb{R}^{n \times d}$  ماتریسی هست که هر سطر آن یک مشاهده از  $d$  ویژگی می باشد و در کل  $n$  مشاهده داریم.  $\mathbf{y} \in \mathbb{R}^n$  بردار برچسب ما می باشد. حال فرض کنید که  $\mathbf{w} \in \mathbb{R}^d$  بردار وزن مدل رگرسیون ما می باشد و  $w^*$  وزن های بهینه می باشند. همچنین فرض کنید مشاهدات ما دارای خاصیت  $\mathbf{X}^T \mathbf{X} = \mathbf{I}$  می باشند. (سفید شده اند) در رگرسیون لاسو بردار وزن های بهینه به صورت زیر به دست میاید :

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} J_{\lambda}(\mathbf{w})$$

$$J_{\lambda} = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

(آ) ابتدا نشان می دهیم که سفیدسازی داده ها باعث مستقل شدن ویژگی ها می شود، به طوری که  $w_i^*$  تنها از ویژگی  $i$ ام نتیجه گیری می شود. برای اثبات این موضوع، ابتدا نشان دهید که  $J_{\lambda}$  را می توان به صورت زیر نوشت:

$$J_{\lambda}(\mathbf{w}) = g(\mathbf{y}) + \sum_{i=1}^d f(\mathbf{X}_{:,i}, \mathbf{y}, w_i, \lambda)$$

که در آن  $\mathbf{X}_{:,i}$  ستون  $i$ ام ماتریس  $\mathbf{X}$  است.

(ب) اگر  $w_i \geq 0$  باشد، مقدار  $w_i$  را بیابید.

(ج) اگر  $w_i < 0$  باشد، مقدار  $w_i$  را بیابید.

(د) با توجه به بخش های قبلی، تحت چه شرایطی  $w_i$  برابر صفر خواهد شد؟ این شرایط چگونه قابل اعمال هستند؟

(ه) همان طور که می دانیم، در رگرسیون ريج، جمله‌ی منظم‌سازی در تابع هزینه به صورت  $\frac{1}{2} \lambda \|\mathbf{w}\|_2^2$  ظاهر می شود. در این حالت،  $w_i$  تحت چه شرایطی برابر صفر می شود؟ تفاوت این حالت با حالت قبلی چیست؟

حل.

(آ) بر اساس روابط داده شده در سوال داریم که :

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} J_{\lambda}(\mathbf{w}), \quad J_{\lambda}(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \cdot \|\mathbf{w}\|_1 \quad (\lambda > 0)$$

$$J_{\lambda}(\mathbf{w}) = \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \cdot \|\mathbf{w}\|_1$$

$$\Rightarrow \frac{1}{2} (\mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}) + \lambda \cdot \|\mathbf{w}\|_1$$

$$J_{\lambda}(\mathbf{w}) = \frac{1}{2} (\mathbf{y}^T \mathbf{y} + (\mathbf{X}\mathbf{w})^T (\mathbf{X}\mathbf{w}) - 2(\mathbf{X}\mathbf{w})^T \mathbf{y}) + \lambda \cdot \|\mathbf{w}\|_1$$

به خاطر سفید سازی ای که انجام داده بودیم داریم که :

$$\mathbf{X}^T \mathbf{X} = \mathbf{I} \Rightarrow J_\lambda(\mathbf{w}) = \frac{1}{2} \mathbf{y}^T \mathbf{y} + \frac{1}{2} \mathbf{w}^T \mathbf{w} - (\mathbf{X} \mathbf{w})^T \mathbf{y} + \lambda \cdot \|\mathbf{w}\|_1$$

$$\mathbf{X} \mathbf{w} = \sum_{i=1}^d \mathbf{X}_i \cdot w_i$$

که در رابطه بالا  $\mathbf{X}_i$  ستون های ماتریس  $\mathbf{X}$  هستند. ضرب  $\mathbf{X} \mathbf{w}$  به ما برداری  $\mathbb{R}^n$  میدهد که هر سطر آن مجموع وزن دار هر سطر ماتریس  $\mathbf{X}$  هست.

$$\mathbf{X} \mathbf{w} = \begin{bmatrix} X_{11} \cdot w_1 + X_{12} \cdot w_2 + X_{13} \cdot w_3 + \dots + X_{1d} \cdot w_d \\ X_{21} \cdot w_1 + X_{22} \cdot w_2 + X_{23} \cdot w_3 + \dots + X_{2d} \cdot w_d \\ X_{31} \cdot w_1 + X_{32} \cdot w_2 + X_{33} \cdot w_3 + \dots + X_{3d} \cdot w_d \\ \dots \\ \dots \\ X_{n1} \cdot w_1 + X_{n2} \cdot w_2 + X_{n3} \cdot w_3 + \dots + X_{nd} \cdot w_d \end{bmatrix}_n = \mathbf{Q}$$

$$(\mathbf{X} \mathbf{w})^T \mathbf{y} = \mathbf{Q}^T \mathbf{y} = \sum_{i=1}^n Q_i \cdot y_i$$

$$\mathbf{w}^T \mathbf{w} = \sum_{i=1}^d w_i^2$$

$$J_\lambda(\mathbf{w}) = \underbrace{\frac{1}{2} \mathbf{y}^T \mathbf{y}}_{g(\mathbf{y})} + \frac{1}{2} \sum_{i=1}^d w_i^2 - \sum_{i=1}^n (\mathbf{X} \mathbf{w})_i^T y_i + \lambda \cdot \sum_{i=1}^d |w_i|$$

برای برطرف کردن اندیس جمع  $\sum_{i=1}^n (\mathbf{X} \mathbf{w})_i^T y_i$  داریم که :

$$(\mathbf{X} \mathbf{w})^T \mathbf{y} = \sum_{i=1}^d w_i \mathbf{X}^T \mathbf{y}$$

$$J_\lambda(\mathbf{w}) = \underbrace{\frac{1}{2} \mathbf{y}^T \mathbf{y}}_{g(\mathbf{y})} + \frac{1}{2} \sum_{i=1}^d w_i^2 - \sum_{i=1}^d w_i \mathbf{X}^T \mathbf{y} + \lambda \cdot \sum_{i=1}^d |w_i|$$

$$\Rightarrow g(\mathbf{y}) + \sum_{i=1}^d \frac{1}{2} \cdot w_i^2 - w_i \cdot \mathbf{X}^T \mathbf{y} + \lambda \cdot |w_i| = g(\mathbf{y}) + \sum_{i=1}^d f(\mathbf{X}_{:,i}, \mathbf{y}, w_i, \lambda)$$

حال برای اینکه نشون دهیم  $w_i^*$  تنها از ویژگی  $i$  نتیجه میشود داریم که :

$$J_\lambda(\mathbf{w}) = \frac{1}{2} \mathbf{y}^T \mathbf{y} + \sum_{i=1}^d \frac{1}{2} w_i^2 - w_i \cdot \mathbf{X}_i^T \mathbf{y} + \lambda \cdot |w_i|$$

$$\frac{\partial J}{\partial w_i} = w_i - \mathbf{X}_i^T \mathbf{y} + \lambda \operatorname{sgn}(w_i)$$

$$\frac{\partial J}{\partial w_i} = 0 \Rightarrow w_i = \mathbf{X}_i^T \mathbf{y} - \lambda \cdot \operatorname{sgn}(w_i)$$

همانطور که در رابطه آخر میبینیم  $w_i$  تنها به ویژگی  $i$ ام بستگی دارد. همچنین میبینیم که فرم بسته ای برای جواب وجود ندارد.

(ب)

$$w_i^* \geq 0 \Rightarrow w_i^* = \mathbf{X}_i^T \mathbf{y} - \lambda$$

(ج)

$$w_i^* < 0 \Rightarrow w_i^* = \mathbf{X}_i^T \mathbf{y} + \lambda$$

(د)

$$w_i^* = 0 \Rightarrow \mathbf{X}_i^T \mathbf{y} = \text{sgn}(w_i^*)\lambda \Rightarrow \lambda = -\text{sgn}(w_i^*)\mathbf{X}_i^T \mathbf{y}$$

اگر مقدار  $\lambda$  را برابر مقدار یافت شده در قسمت قبل قرار دهیم،  $w_i^*$  صفر میشود.

(ه)

$$\begin{aligned} R &= \frac{1}{2}\lambda\|\mathbf{w}\|_1^2 \Rightarrow J_\lambda(\mathbf{w}) = \frac{1}{2}\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \frac{1}{2}\lambda\|\mathbf{w}\|_1^2 \\ J_\lambda &= \frac{1}{2}(\mathbf{y}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2(\mathbf{X}\mathbf{w})^T \mathbf{y}) + \frac{1}{2}\mathbf{w}^T \mathbf{w} \\ &= \frac{1}{2} \sum_{j=1}^n y_j^2 + \frac{1}{2} \sum_{i=1}^d w_i^2 - \sum_{i=1}^d w_i \cdot (\mathbf{X}_{:,i}^T \mathbf{y}) + \frac{1}{2}\lambda \sum_{i=1}^d w_i^2 \\ \frac{\partial J}{\partial w_i} &= 0 \Rightarrow w_i - \mathbf{X}_{:,i}^T \mathbf{y} + \lambda \cdot w_i = 0 \Rightarrow w_i = \frac{\mathbf{X}_{:,i}^T \mathbf{y}}{1 + \lambda} \end{aligned}$$

در مسیله ridge تنها وقتی که برچسب داده ها بر ویژگی داده عمود باشد،  $w_i^* = 0$  میشود ( $\mathbf{X}_{:,i}^T \mathbf{y} = 0$ ) ولی در مسیله لاسو میتوانیم با انتخاب بهینه  $\lambda$  که در بخش قبل به دست آوردیم به این نتیجه برسیم.

۷. (۱۵ نمره) تابع زیر را در نظر بگیرید :

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2 x_2)}}$$

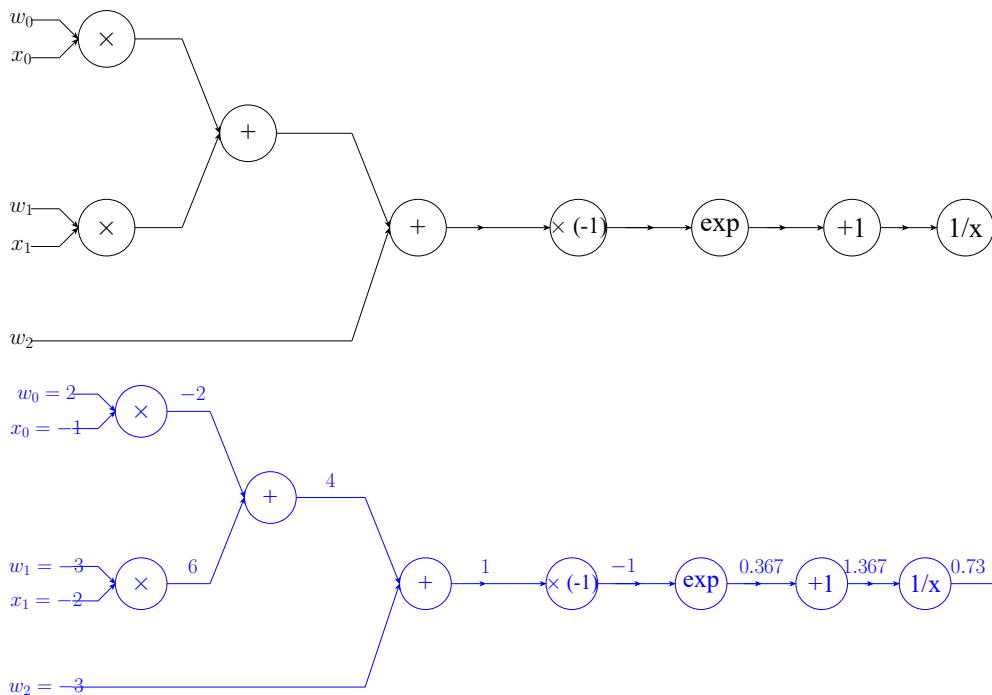
شبکه متناظر با این تابع در زیر آمده است. با فرض ورودی های

$$w_0 = 2, x_0 = -1, w_1 = -3, x_1 = -2, w_2 = -3$$

ابتدا خروجی نهایی را محاسبه کنید. سپس با استفاده از Backpropagation مشتق خروجی نهایی را نسبت به خروجی هر راس حساب کنید.

حل. در هر مرحله از forward-pass ما بر اساس گراف محاسباتی مان خروجی را محاسبه میکنیم که به فرم زیر در میاید :

حال در هر مرحله از Backpropagation نسبت به تابع گره ای که در آن هستیم مشتق میگیریم و به عقب حرکت میکنیم. (برای  $w$  را مینویسیم)



$$g(x) = \frac{1}{x}, h(x) = \exp(ax), Q(x) = ax, m(x) = b + x$$

$$\frac{\partial \hat{y}}{\partial w_0} = \frac{\partial \hat{y}}{\partial g} \cdot \frac{\partial g}{\partial w_0} = \frac{\partial \hat{y}}{\partial g} \cdot \frac{\partial g}{\partial m} \cdot \frac{\partial m}{\partial w_0} = \frac{\partial \hat{y}}{\partial g} \cdot \frac{\partial g}{\partial m} \cdot \frac{\partial m}{\partial h} \cdot \frac{\partial h}{\partial w_0} = \frac{\partial \hat{y}}{\partial g} \cdot \frac{\partial g}{\partial m} \cdot \frac{\partial m}{\partial h} \cdot \frac{\partial h}{\partial Q} \cdot \frac{\partial Q}{\partial w_0}$$

$$1) \frac{\partial \hat{y}}{\partial w_0} = \frac{\partial \hat{y}}{\partial g} \Big|_{g=1.367} \cdot \frac{\partial g}{\partial w_0} = \frac{-1}{x^2} \Big|_{x=1.367} \cdot \frac{\partial g}{\partial w_0} = -0.53 \cdot \frac{\partial g}{\partial w_0}$$

$$2) \frac{\partial \hat{y}}{\partial g} \cdot \frac{\partial g}{\partial m} \Big|_{m=0.367} \cdot \frac{\partial m}{\partial w_0} = -0.53 \cdot 1 \cdot \frac{\partial m}{\partial w_0}$$

$$3) \frac{\partial \hat{y}}{\partial g} \cdot \frac{\partial g}{\partial m} \cdot \frac{\partial m}{\partial h} \cdot \frac{\partial h}{\partial w_0} = -0.53 \cdot 1 \cdot \exp(x) \Big|_{-1} \cdot \frac{\partial h}{\partial w_0} = -0.53 \cdot 1 \cdot e^{-1} = -0.194 \frac{\partial h}{\partial w_0}$$

$$4) \frac{\partial \hat{y}}{\partial g} \cdot \frac{\partial g}{\partial m} \cdot \frac{\partial m}{\partial h} \cdot \frac{\partial h}{\partial Q} \cdot \frac{\partial Q}{\partial w_0} = -0.194 \cdot \frac{\partial Q}{\partial w_0} \Big|_{a=-1} = -0.194 \cdot -1 = 0.194 = T$$

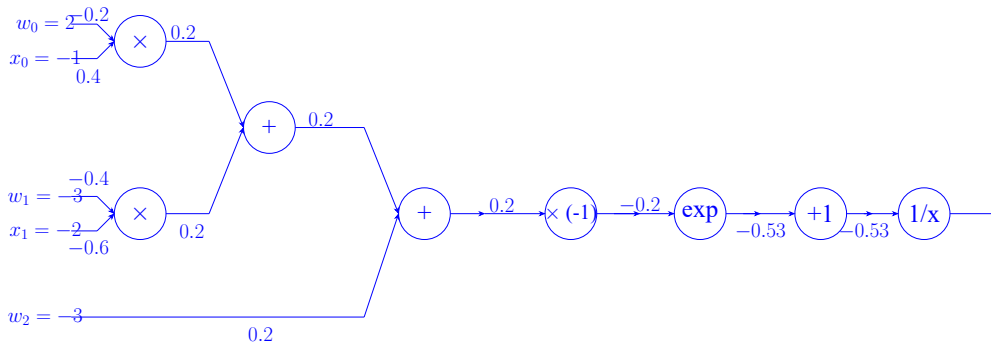
$$5) \frac{\partial T}{\partial m} \cdot \frac{\partial m}{\partial w_0} = 0.194$$

$$6) \frac{\partial T}{\partial m} \cdot \frac{\partial m}{\partial w_0} = 0.194$$

$$7) \frac{\partial T}{\partial m} \cdot \frac{\partial m}{\partial w_0} = 0.194 \cdot x_0 \Big|_{x_0 = -1} = -0.194$$

شماره گذاری برای راحت تر شدن خواندن مراحل انجام شده است و صرفاً برای نمایش یک دور backpropagation می باشند و نیازی به نوشتن آنها در پاسخ نبوده است و درستی جواب آخر کافی می باشد. همانطور که دیدید ۷ مشتق جزیی نسبت به توابع و اجزای مختلف شبکه در مسیر متغیر  $w$  گرفتیم تا نرخ تغییرات با توجه به متغیر مشخص شده را بیابیم که این تعداد برابر تعداد نود های بین متغیر مورد نظر و خروجی شبکه بود. مقادیر نهایی Backpropagation در هر مرحله در زیر آمده است : (مقادیر به دو رقم اعشار گرد شده اند)

با استفاده از مقادیر به دست آمده میتوان متغیر ها را طوری تغییر داد تا به جواب بهینه مان برسیم (الگوریتم گرادیان کاهشی نمونه ای از این دسته روش ها می باشد).



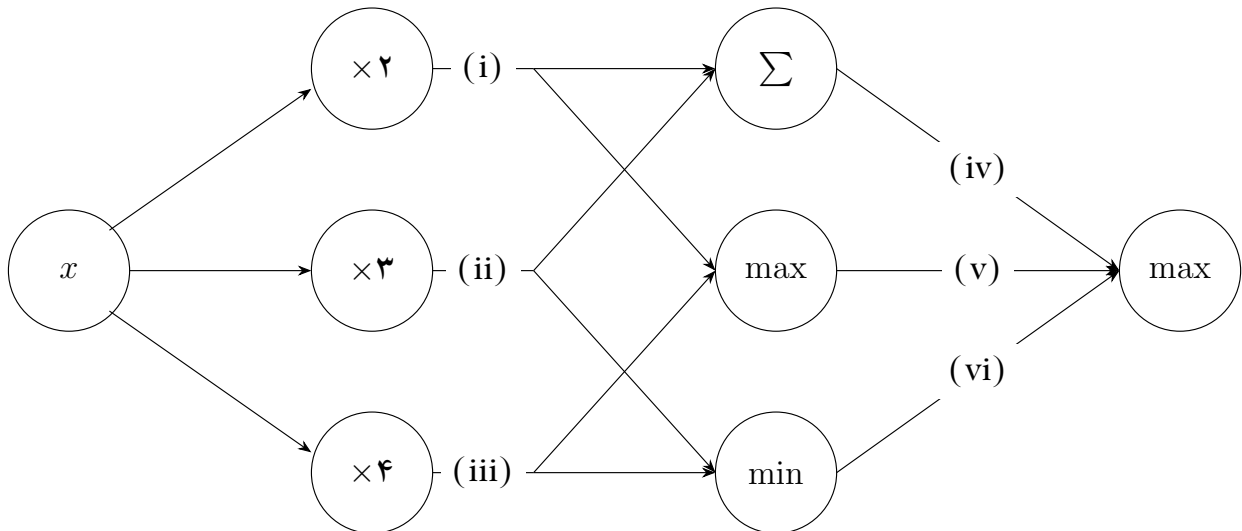
۸. (۱۵+۵ نمره)

(آ) عملیات Forward Propagation را برای شبکه‌ی عصبی زیر با ورودی  $x = ۱$  انجام دهید.

(۱) مقادیر مربوط به (i), (ii), (iii) را به دست آورید.

(۲) مقادیر مربوط به (iv), (v), (vi) را به دست آورید.

(۳) مقدار خروجی مدل را به دست آورید.



(ب) (امتیازی) شکل زیر یک شبکه‌ی عصبی با وزن‌های  $a, b, c, d, e, f$  را نشان می‌دهد. ورودی این شبکه  $x_1$  و  $x_2$  هستند. لایه‌ی مخفی نخست به شکل زیر محاسبه می‌شود.

$$r_1 = \max(c \cdot x_1 + e \cdot x_2, *) \quad r_2 = \max(d \cdot x_1 + f \cdot x_2, *)$$

خروجی لایه‌ی مخفی دوم نیز به شکل زیر است.

$$s_1 = \frac{1}{1 + \exp(-a \cdot r_1)} \quad s_2 = \frac{1}{1 + \exp(-b \cdot r_2)}$$

خروجی شبکه نیز برابر  $y = s_1 + s_2$  است.

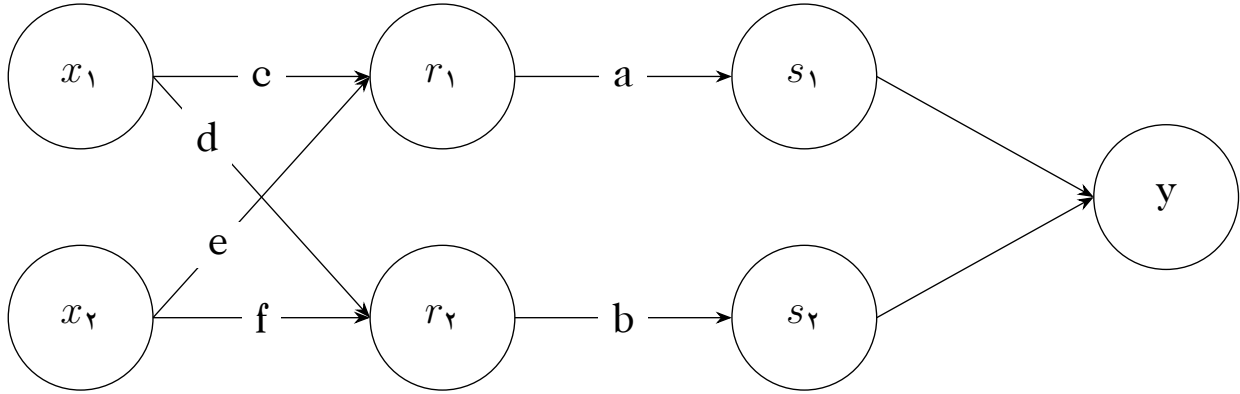
فرض کنید ورودی شبکه  $x_1 = ۱, x_2 = -۱$  و مقدار وزن‌ها برابر

$$a = ۱, b = ۱, c = ۴, d = ۱, e = ۲, f = ۲$$

باشد. در این صورت خروجی تقریبی شبکه به صورت

$$r_1 = 2, r_2 = 0, s_1 = 0.9, s_2 = 0.5, y = 1/4$$

خواهد بود.



با استفاده از مقادیر تقریبی داده شده از الگوریتم backpropagation استفاده کنید تا مقدار مشتق‌های جزئی را حساب کنید. مقادیر را به صورت عددی به دست آورید.

(۱) مقدار  $\frac{\partial y}{\partial b}$  و  $\frac{\partial y}{\partial a}$  را به دست آورید.

(۲) مقدار  $\frac{\partial r_1}{\partial e}$  و  $\frac{\partial r_1}{\partial c}$  را به دست آورید.

(۳) مقدار  $\frac{\partial r_2}{\partial f}$  و  $\frac{\partial r_2}{\partial d}$  را به دست آورید.

(۴) مقدار  $\frac{\partial y}{\partial r_2}$  و  $\frac{\partial y}{\partial r_1}$  را به دست آورید.

(۵) با استفاده از مقادیر به دست آمده در بخش ۴ و خاصیت زنجیره‌ای مشتقات،  $\frac{\partial y}{\partial c}$  و  $\frac{\partial y}{\partial e}$  را به دست آورید.

(۶) با استفاده از مقادیر به دست آمده در بخش ۴ و خاصیت زنجیره‌ای مشتقات،  $\frac{\partial y}{\partial d}$  و  $\frac{\partial y}{\partial f}$  را به دست آورید.

(۷) تمام وزن‌های شبکه را یک مرحله با نرخ یادگیری یک بروزرسانی کنید.

حل.

(۱) (I)

$$(i) : x \times 2 = 2, \quad (ii) : x \times 3 = 3, \quad (iii) : x \times 4 = 4$$

(۲)

$$(iv) : (i) + (ii) = 5, \quad (v) : \max((i), (iii)) = 4, \quad (vi) : \min((ii), (iii)) = 3$$

(۳)

$$\text{out} = \max((iv), (v), (vi)) = 5$$

(۱) (ب)

$$\frac{\partial y}{\partial a} = \frac{\partial y}{\partial s_1} \frac{\partial s_1}{\partial a} = 1 \cdot r_1 s_1 (1 - s_1) = 0.18$$

$$\frac{\partial y}{\partial b} = \frac{\partial y}{\partial s_2} \frac{\partial s_2}{\partial b} = 1 \cdot r_2 s_2 (1 - s_2) = 0$$

(2)

$$c \cdot x_{\text{I}} + e \cdot x_{\text{V}} > 0 \Rightarrow r_{\text{I}} = c \cdot x_{\text{I}} + e \cdot x_{\text{V}} \Rightarrow \frac{\partial r_{\text{I}}}{\partial c} = x_{\text{I}} = 1, \quad \frac{\partial r_{\text{I}}}{\partial e} = x_{\text{V}} = -1$$

(3)

$$d \cdot x_{\text{I}} + f \cdot x_{\text{V}} = -1 < 0 \Rightarrow \frac{\partial r_{\text{V}}}{\partial d} = \frac{\partial r_{\text{V}}}{\partial f} = 0$$

(4)

$$\frac{\partial y}{\partial r_{\text{I}}} = \frac{\partial y}{\partial s_{\text{I}}} \frac{\partial s_{\text{I}}}{\partial r_{\text{I}}} = 1 \cdot a s_{\text{I}} (1 - s_{\text{I}}) = 0.09$$

$$\frac{\partial y}{\partial r_{\text{V}}} = \frac{\partial y}{\partial s_{\text{V}}} \frac{\partial s_{\text{V}}}{\partial r_{\text{V}}} = 1 \cdot b s_{\text{V}} (1 - s_{\text{V}}) = 0.25$$

(5)

$$\frac{\partial y}{\partial c} = \frac{\partial y}{\partial r_{\text{I}}} \frac{\partial r_{\text{I}}}{\partial c} = 0.09 \times 1 = 0.09$$

$$\frac{\partial y}{\partial e} = \frac{\partial y}{\partial r_{\text{I}}} \frac{\partial r_{\text{I}}}{\partial e} = 0.09 \times -1 = -0.09$$

(6)

$$\frac{\partial y}{\partial f} = \frac{\partial y}{\partial r_{\text{V}}} \frac{\partial r_{\text{V}}}{\partial f} = 0$$

$$\frac{\partial y}{\partial d} = \frac{\partial y}{\partial r_{\text{V}}} \frac{\partial r_{\text{V}}}{\partial d} = 0$$

(7)

$$a \leftarrow 0.82, b \leftarrow 1, c \leftarrow 3.91, d \leftarrow 1, e \leftarrow 0.09, f \leftarrow 2$$