

تمرین سری چهارم (رس) داده کاوی - فاطمه غلامزاده - ۹۵۳۱۰۶۰

### سوال ۱

تقنیه هم‌گرایی: اگر مجموعه وزن‌های  $w^*$  جدایی پذیر خطی باشند (یعنی قابلیت جداسازی یک مجموعه‌ی محدود را داشته باشند)، الگوریتم یادگیری پرسپترون همگرا خواهد شد.

تقنیه: اگر بردار وزن  $w^*$  را به گونه‌ای داشته باشیم به ازای حد  $p$ :

$$f(x(p) \cdot w^*) = t(p)$$

آن‌گاه برای هر بردار اولیه  $w$  الگوریتم یادگیری پرسپترون به بردار وزنی نزدیک می‌شود که برای همه‌ی الگوها پاسخ صحیح می‌دهد (در مرحله‌ی مناسبت).

$p$ : مقدار هدف معادل بردارها  
 $x(p)$ : بردارهای ورودی  
 $f$ : activation function، خروجی

نکته: وزن‌هایی که الگوریتم به آن‌ها نزدیک می‌شود لزوماً محفوز فرد و لزوماً برابر با  $w^*$  نیستند.

$$\text{prediction}(y') = \begin{cases} 1 & wx + b \geq 0 \\ 0 & wx + b < 0 \end{cases}$$

### سوال ۲

NOR:  $A \Rightarrow B \Rightarrow \overline{A+B}$

گیت NOR فقط در صورتی که هر دو ورودی ۰ باشند ۱ می‌دهد.

A	B	$\overline{A+B}$
0	0	1
0	1	0
1	0	0
1	1	0

سطر اول: مقدار اولیه  
 $w_1 = 1$   
 $w_2 = 1$   
 $b = -1$

$$\text{model} = x_1(1) + x_2(1) - 1$$

سطر اول:  $x_1 = 0 \Rightarrow 0 + 0 - 1 = -1 \Rightarrow wx + b = -1 < 0 \Rightarrow y' = 0$



با توجه به فرضی جدول کمی خواهیم داشت  $x$  پس  $x$  را تغییر

می دهیم در برابر  $1$  در نظریه می گیریم:  $\checkmark y' = 1 \Rightarrow 0 + 0 + 1 = 1 > 0$

سطر دوم:  $x_1 = 0$  و  $x_2 = 1 \Rightarrow wx + b = 0 + 1 + 1 = 2 > 0 \Rightarrow y' = 1$

طبق ردیف دوم فرضی باید صغری نشد پس خط  $x$  این بار  $w_1$  و  $w_2$  را

به  $-1$  تغییر می دهیم و  $b = 0.5$   $wx + b = 0 - 1 + 0.5 = -0.5 < 0 \Rightarrow y' = 0 \checkmark$

سطر سوم:  $x_1 = 1$   $x_2 = 0$

$wx + b = -1 + 0 + 0.5 = -0.5 < 0 \Rightarrow y' = 0 \checkmark$

با توجه به ردیف سوم که صغری است جواب درستی حاصل شد

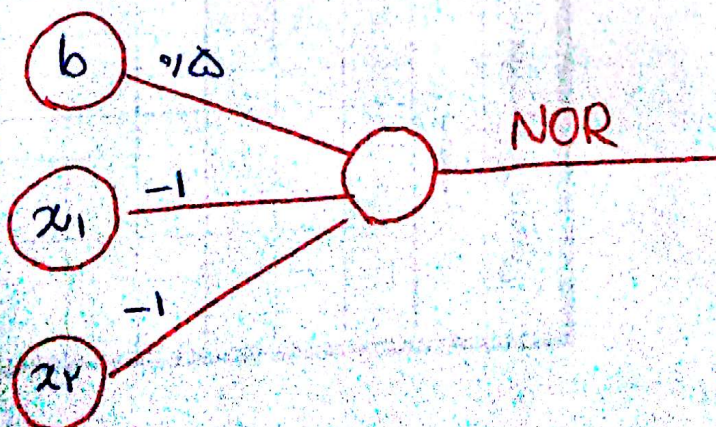
سطر چهارم:  $x_1 = 1$   $x_2 = 1$

$wx + b = -1 - 1 + 0.5 = -1.5 < 0 \rightarrow y' = 0 \checkmark$

NOR Gate Model:

$$\begin{cases} w_1 = w_2 = -1 \\ b = 0.5 \end{cases}$$

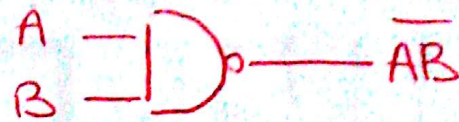
$$y' = \begin{cases} 1 & -x_1 - x_2 + 0.5 > 0 \\ 0 & -x_1 - x_2 + 0.5 < 0 \end{cases}$$





NAND:

A	B	out
0	0	1
0	1	1
1	0	1
1	1	0



سطر اول:

$$\begin{cases} w_1 = w_2 = 1 \\ b = -1 \end{cases}$$

مقداردهی اولیه  $w_1, w_2$  و  $b$ :

$$x_1 = x_2 = 0 \rightarrow wx + b = 0 + 0 + (-1) = (-1) < 0 \rightarrow y' = 0$$

غلط

$$wx + b = 0 + 0 + 1 = 1 > 0 \rightarrow y' = 1 \checkmark$$

با آیدیت می کنیم  $b = 1$

$$x_1 = 0, x_2 = 1 \rightarrow wx + b = 0 + 1 + 1 = 2 > 0 \rightarrow y' = 1 \checkmark$$

سطر دوم:

$$x_1 = 1, x_2 = 0 \Rightarrow wx + b = 1 + 0 + 1 = 2 > 0 \rightarrow y' = 1 \checkmark$$

سطر سوم:

$$x_1 = x_2 = 1$$

$$wx + b = 1 + 1 + 1 = 3 > 0 \rightarrow y' = 1 \quad X$$

$$w_1 = w_2 = -1$$

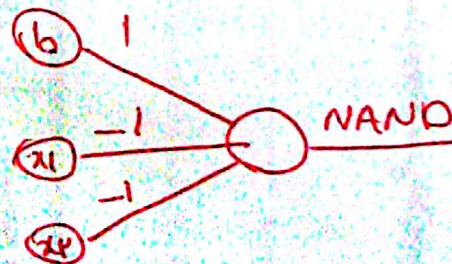
$w_1$  و  $w_2$  را آیدیت می کنیم

$$wx + b = -1 - 1 + 1 = -1 < 0 \rightarrow y' = 0 \checkmark$$

NAND Gate Model:

$$w_1 = w_2 = -1$$

$$b = 1$$

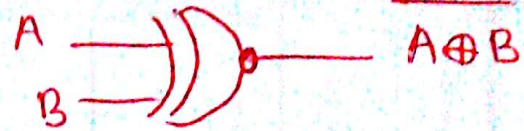


$$y' = \begin{cases} 1 & -x_1 - x_2 + 1 > 0 \\ 0 & -x_1 - x_2 + 1 < 0 \end{cases}$$



XNOR :

A	B	out
0	0	1
0	1	0
1	0	0
1	1	1



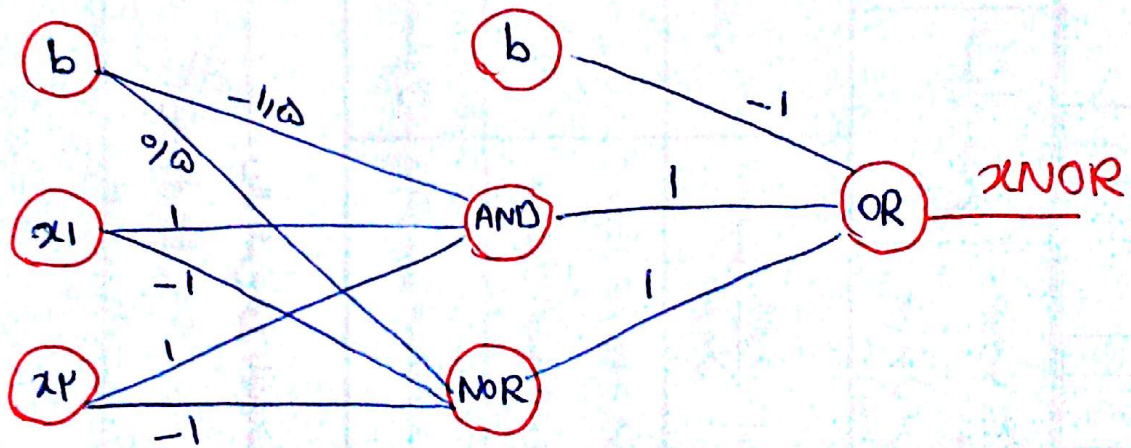
$$\overline{A \oplus B} = \underbrace{AB}_{\text{AND}} + \underbrace{A'B'}_{\text{NOR}}$$

گیت xnor از گیت های AND و NOR و OR تشکیل شده :  
پس باید پرسترون های آن ها را  
ما هم ترکیب کنیم.

$$\text{AND} \rightarrow x_1 + x_2 - 1/5$$

$$\text{NOR} \rightarrow -x_1 - x_2 + 0/5$$

$$\text{OR} \rightarrow x_1 + x_2 - 1$$



سوال ۳ | الگوریتم Kmean++ : تفاوت این الگوریتم با K-means

نقطه در انتخاب مختصات اولیه مراکز خوشه ها است. ابتدا یکی از داده ها به طور تصادفی انتخاب شده و به عنوان مرکز یک خوشه در نظر گرفته می شود. سپس مراکز بعدی با احتمالی که با مربع فاصله آن داده با نزدیکترین مرکز خوشه از بین مراکز که تا به حال انتخاب شده اند، متناسب است، انتخاب می شود و این روند ادامه می یابد تا K تا مرکز انتخاب شود. این الگوریتم را ابتدا زمان بیشتری را صرف انتخاب مراکز می کند (نسبت به Kmean) اما این کار بعداً باعث می شود الگوریتم زودتر همگرا شود و همیشه سرعت بیشتری از K-means دارد و خطای کمتری هم دارد. زمان اجرا:  $O(\log(K))$



**سوال ۱۴** در ایجاد بالا ۲ روش برای بررسی سازی اطلاعات داریم :

PCA ① T-SNE ②

① PCA → در این روش دیتا را به فضای دیگری می‌بریم و Component ها را بر اساس اهمیتشان Sort کرده و Component های اصلی را افزایش می‌کنیم به گونه‌ای که بخش عمده‌ای از دیتای اصلی توسط این Component های اصلی قابل بازسازی است. این روش برای بررسی سازی کلاس‌های فضای اصلی ۲ و ۳ بعدی به کار می‌کند. این روش بیشتر برای اطلاعات تا ۳ بعد مناسب است و در ایجاد بالا خوب عمل نمی‌کند.

② T-SNE : برای داده‌های با ابعاد بالا مناسب است. از ارتباط مکانی بین نقاط برای نگاشت به ابعاد پایین‌تر استفاده می‌کند و ساختاری غیرخطی دارد. این روش یک توزیع احتمال (مثلاً گوسی) ایجاد می‌کند که نشان دهنده‌ی ارتباط میان نقاط همسایه است، سپس یک فضای با ابعاد کم را می‌سازد که آن توزیع را تا جایی که می‌تواند حفظ کند و از دست دادن استفاده می‌کند.

یک روش دیگر برای نمایش اطلاعات در ابعاد بالا، Correlation heatmap نام دارد. ما رسم کردن همبستگی‌ها را برای هر فیچر می‌توان ایده گرفت که هر بعد چگونه توزیع شده است.

Correlation heatmap نشان دهنده‌ی چگونگی ارتباط دو بعد از داده‌ها با هم است. اگر ضریب همبستگی ۱ بعد زیاد باشد نشان می‌دهد که ویژگی‌ها یکدیگر وابسته‌اند اما اگر کم باشد استقلال نسبی آن‌ها را نتیجه می‌دهد.

برای نشان دادن داده‌ها در ۳ بهای برای ویژگی سوم می‌توانیم از شکل،  
رنگ و ... استفاده کنیم.