

## هوش محاسباتی - ترینیداد

1.  $T(u, v) = (1 + u, v)$

$$\left. \begin{aligned} T((u_1, y_1), (u_2, y_2)) &= T(u_1 + u_2, y_1 + y_2) = (1 + u_1 + u_2, y_1 + y_2) \\ T(u_1, y_1) + T(u_2, y_2) &= (1 + u_1, y_1) + (1 + u_2, y_2) = (2 + u_1 + u_2, y_1 + y_2) \end{aligned} \right\}$$

جمع پذیر نیست  $\rightarrow$  عمل غیر خطی

•  $T(u_1, u_2) = (u_2, u_1)$

$$\left. \begin{aligned} T((u_1, y_1) + (u_2, y_2)) &= T(u_1 + u_2, y_1 + y_2) = (y_1 + y_2, u_1 + u_2) \\ T(u_1, y_1) + T(u_2, y_2) &= (y_1, u_1) + (y_2, u_2) = (y_1 + y_2, u_1 + u_2) \end{aligned} \right\} \text{جمع پذیر}$$

$$\left. \begin{aligned} T(c \cdot (u, y)) &= T(cu, cy) = (cy, cu) \\ c \cdot T(u, y) &= c \cdot (y, u) = (cy, cu) \end{aligned} \right\} \text{تغییر در نتیجه خطی است}$$

•  $T(u, u) = (u^2, u)$

$$\left. \begin{aligned} T((u_1, y_1) + (u_2, y_2)) &= T(u_1 + u_2, y_1 + y_2) = ((u_1 + u_2)^2, y_1 + y_2) \\ T(u_1, y_1) + T(u_2, y_2) &= (u_1^2, y_1) + (u_2^2, y_2) = (u_1^2 + u_2^2, y_1 + y_2) \end{aligned} \right\} \begin{array}{l} \text{جمع ناپذیر} \\ \text{غیر خطی} \end{array}$$

•  $T(u, u) = (\sin u, u)$

$$\left. \begin{aligned} T((u_1, y_1) + (u_2, y_2)) &= T(u_1 + u_2, y_1 + y_2) = (\sin(u_1 + u_2), y_1 + y_2) \\ T(u_1, y_1) + T(u_2, y_2) &= (\sin u_1, y_1) + (\sin u_2, y_2) = (\sin u_1 + \sin u_2, y_1 + y_2) \end{aligned} \right\} \begin{array}{l} \text{جمع ناپذیر} \\ \text{غیر خطی} \end{array}$$

•  $T(u, u) = (u, -u, 0)$  مشابه قسمت دوم

$$T((u_1, y_1) + (u_2, y_2)) = (u_1 + u_2, -y_1 - y_2, 0) = T(u_1, y_1) + T(u_2, y_2)$$

$$T(cu, cy) = (cu, -cy, 0) = c(u, -y, 0) = c \cdot T(u, y)$$

2. می توانیم دو ماتریس را به فرم RREF تبدیل کنیم. اگر RREF آن ها یکسان باشد هم ارزشی دارند.

RREF(A):

$$A = \begin{bmatrix} 2 & 0 & 0 \\ a & -1 & 0 \\ b & c & 3 \end{bmatrix} \xrightarrow{r_1 \div 2} \begin{bmatrix} 1 & 0 & 0 \\ a & -1 & 0 \\ b & c & 3 \end{bmatrix} \xrightarrow{\substack{r_2 - ar_1 \\ r_3 - br_1}} \begin{bmatrix} 1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & c & 3 \end{bmatrix} \xrightarrow{\substack{r_2 \times -1 \\ r_3 - cr_2}} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

RREF(B):

$$B = \begin{bmatrix} 1 & 1 & 2 \\ 0 & 2 & 3 \\ 0 & 2 & 3 \end{bmatrix} \xrightarrow{r_2 \div 2} \begin{bmatrix} 1 & 1 & 2 \\ 0 & 1 & 1.5 \\ 0 & 2 & 3 \end{bmatrix} \xrightarrow{r_3 - 2r_2} \begin{bmatrix} 1 & 1 & 2 \\ 0 & 1 & 1.5 \\ 0 & 0 & 0 \end{bmatrix} \xrightarrow{r_1 - r_2} \begin{bmatrix} 1 & 0 & 0.5 \\ 0 & 1 & 1.5 \\ 0 & 0 & 0 \end{bmatrix}$$

$\Rightarrow RREF(A) \neq RREF(B) \Rightarrow$  هم ارزش سطر نیستند

$$A = PDP^{-1}, \quad P = \begin{bmatrix} 1 & 0 & \frac{1}{2} \\ 0 & -1 & 1 \\ 1 & 1 & 0 \end{bmatrix}, \quad D = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix} \quad .3$$

$$P^{-1} = \frac{1}{\det(P)} \cdot \text{adj}(P) \quad \Rightarrow \quad \det(P) = 1 \times (-1) \times 0 + 0 \times 1 \times 1 + \frac{1}{2} \times 0 \times 1 - (\frac{1}{2} \times (-1) \times 1) - (1 \times 1 \times 1) - (0) = \frac{1}{2}$$

$$\text{adj}(P) = \begin{bmatrix} -1 & \frac{1}{2} & \frac{1}{2} \\ 1 & -\frac{1}{2} & \frac{1}{2} \\ 1 & 1 & -1 \end{bmatrix} \quad \Rightarrow \quad P^{-1} = \frac{1}{\frac{1}{2}} \times \text{adj}(P) = \begin{bmatrix} -2 & 1 & 1 \\ 2 & -1 & -1 \\ 2 & 2 & -2 \end{bmatrix}$$

$$\Rightarrow A = PDP^{-1} = \begin{bmatrix} 1 & 0 & \frac{1}{2} \\ 0 & -1 & 1 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix} \begin{bmatrix} -2 & 1 & 1 \\ 2 & -1 & -1 \\ 2 & 2 & -2 \end{bmatrix} = \begin{bmatrix} 1 & 4 & -2 \\ 2 & 8 & -4 \\ 2 & -1 & -1 \end{bmatrix}$$

$$\Rightarrow A \begin{bmatrix} 3 \\ -1 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 & 4 & -2 \\ 2 & 8 & -4 \\ 2 & -1 & -1 \end{bmatrix} \begin{bmatrix} 3 \\ -1 \\ 2 \end{bmatrix} = \begin{bmatrix} -5 \\ -10 \\ 5 \end{bmatrix}$$

Scanned with



#### Gradient Decent :

4.

مزایا	معایب
+ یاباری بیشتر به علت استفاده از همه داده ها.	- چون روی همه داده ها کار می کند هزینه بالاتری دارد.
+ در توابع محدب مینیمم گلوبال را پیدا می کند.	- حافظه بیشتر نیاز دارد و در مجموعه های بزرگ کند است.
+ نوسان کمتری دارد.	- ممکن است در مینیمم های محلی گیر کند.

#### Stochastic Gradient Decent (batch-size = 1)

مزایا	معایب
+ هر تکرار سریع انجام می شود.	- نوسان زیاد در مسیر.
+ مصرف حافظه کم.	- همگرا نیست غیر قطعی: ممکن است حول نقطه بچرخد.
+ می تواند از مینیمم های محلی فرار کند.	- نیاز به تنظیم نرخ یادگیری.
+ مناسب برای داده های حجیم و یادگیری آنلاین.	- حساسیت به نویز.

در روش GD ماده پارامترهای مدل بروی تمامی داده های آموزشی بروز می شوند ولی در SGD به اندازه batch-size که اینجا یک در نظر گرفتیم.

#### 5. ابتدا دو مورد MSE را بررسی می کنیم. اولین و اصلی ترین مشکل این روش این است که ترکیب

آن با سیگموئید یک تابع غیر محدب است. این به این معناست که الگوریتم بهینه سازی نوی تری نیاز داریم تا در مینیمم های محلی گیر نکند. مشکل دیگری که می توان فهمید این است که

MSE در رگرسیون لاجستیک ناسازگاری آماری دارد. یعنی MSE فرض می کند خطاها

توزیع نرمال دارند در صورتی که خروجی ها 0 و 1 هستند و توزیع برنولی دارند.

از طرفی Maximum Likelihood ترکیبش با سگمویید محذب است و همچنین از لحاظ آماری

سازگاری دارد.

6. باید نقطه ای را پیدا کنیم که اگر حذفش کنیم، شیب خط بیشترین تغییر را بکند.

این به این معناست که هرچه نقطه ای از شیب حاصله بیشتری داشته باشد تاثیر بیشتری

هم دارد. با توجه به تجمع نقاط در نمودار، بدون محاسبه می توان فهمید که دو نقطه ای

بالا سمت راست نمودار بیشترین حاصله را دارند. پس با همین فرض می توانیم حدس بزنیم

که نقطه سمت راست نمودار همان نقطه مدنظر است (حدوداً (9.6, 4.5)). برای اثبات

آن کافیت شیب را در حالت عادی و در حالت بدون این نقطه حساب کنیم و ببینیم که بیشترین

تغییر را دارد.



7. سوال 1: اگر نرخ یادگیری ثابت و بسیار بزرگ باشد ممکن است حول نقطه بهینه نوسان کند.

و اگر کوچک باشد همگرا این به بندگی انجام می شود یا در مینیمم محلی گیر می کند. برای حل این

مشکل می توانیم نرخ یادگیری را با الگوریتم های مختلف تغییر دهیم مانند Adam

سوال 2: این دو بردار برهم عمود هستند. مثال:

$$A = [1, 0], B = [0, 1] \Rightarrow A \cdot B = 1 \times 0 + 0 \times 1 = 0$$

می توان دید که عمود هستند

سوال 3:

نادرست در رگرسیون خطی اگر واریانس خطاها ناهمسان باشد تخمین های مدل

ناایمن یا متعصب می باشد ولی واریانس تخمین ها بهینه نیست و بیش بینی ها قابل اعتماد نیستند.

برای حل این مشکل می توان از روش هایی مانند وزن دهی یا تبدیل متغیر استفاده کرد

سوال 4: این تابع برای مسئله طبقه بندی بهتر است چون تفاوت بین توزیع واقعی و

بیش بینی شده را اندازه می گیرد و به مدل کمک می کند تا احتمال های طیفات را بهتر تخمین بزند.

دلیل هاهنگی این تابع این است که برای خروجی های بین منفرد یک طراحی شده و اگر

بیش بینی مدل با مقدار واقعی تفاوت داشته باشد جریمه سنگینی اعمال می کند.

سوال 5: گاهی بیش از حد نرخ یا دلیری ممکن است باعث کند شدن همگرایی شود یا

حتی در مینیمم های محلی گیر کند. همگرایی به چیزهای دیگری مانند ساختار تابع هزینه و

مقدار اولیه پارامترها هم وابسته است.

8. الف) در الگوریتم خطی هدف یافتن بردار وزن های  $w$  است که خطای پیش بینی را کمینه کند.

مانند به فرض استفاده از  $J$  داریم:  $J = \frac{1}{2} \|y - Xw\|^2$  در نتیجه برای خطای مربعات داریم:  $J = \frac{1}{2} (y - Xw)^T (y - Xw)$

برای کمینه کردن این تابع مشتق  $J$  نسبت به  $w$  را صفر می کنیم.

$$\frac{\partial J}{\partial w} = -X^T (y - Xw) = 0 \Rightarrow X^T y = X^T X w \Rightarrow \frac{X^T y}{X^T X} \rightarrow \frac{X^T y}{X^T X}$$

ب) اگر ویژگی ها مستقل باشد یعنی سطرهای ماتریس  $X$  مستقل خطی هستند آنگاه

ماتریس  $X^T X$  یک ماتریس قطری خواهد بود زیرا ضرب داخلی سطرهای مستقل خطی صفر

است. در این حالت بردار وزن بهینه  $w$  برای الگوریتم با همه ویژگی ها به صورت

$$w = (X^T X)^{-1} X^T y$$

زیر است:

که اگر  $X^T X$  قطری باشد داریم:  $w = \frac{X^T y}{X^T X}$  که همین نتیجه ای است که در الف گرفتیم.



ج. در این مدل  $W$  یک بردار وزن و  $w_0$  بایاس است.

$$J(W, w_0) = \|y - (X^T W + w_0)\|^2$$

این تابع هزینه ما را کمینه می کنیم:

$$\frac{\partial J}{\partial W} = -2X(y - X^T W - w_0)^T = 0, \quad \frac{\partial J}{\partial w_0} = -2 \sum_{i=1}^N (y_i - W^T x_i - w_0) = 0$$

$$\Rightarrow w_0 = \bar{y} - W^T \bar{x}, \quad \Rightarrow W = (XX^T)^{-1} X(y - w_0)$$

9. با فرض  $y = a + bx$  داریم: ابتدا میابیم ما را حساب می کنیم:

$$\bar{x} = \frac{521}{8} = 65.125$$

$$\bar{y} = \frac{232}{8} = 29$$

$$b = \frac{n \sum xy - (\sum x)(\sum y)}{n \sum x^2 - (\sum x)^2}$$

شیب  $b$ :

$$\sum xy = 14971, \quad \sum x^2 = 34203$$

$$\Rightarrow b = \frac{8 \times 14971 - 521 \times 232}{8 \times 34203 - 501^2} = \frac{-1104}{2183} \approx -0.506$$

$$a = \bar{y} - b\bar{x} = 29 - (-0.506 \times 65.125) \approx 61.95$$

عرض از مبدا  $a$ :

$$\Rightarrow y = a + bx = 61.95 - 0.506x$$

10. در روش SGD با فرض اینکه تابع هزینه کل به صورت زیر تعریف شود:

$$\sum_{i=1}^n J_i(\theta) \quad \frac{1}{n} = J(\theta)$$

در روش SGD به جای  $\frac{1}{n} \sum_{i=1}^n \nabla J_i(\theta)$  که گرادینت کل است، ما از یک ست نمونه

استفاده می‌کنیم.  $(\nabla J_i(\theta))$  حالا اگر داده‌ها به صورت تصادفی انتخاب شوند داریم:

$$\nabla J(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla J_i(\theta) = E_i[\nabla J_i(\theta)]$$

در نتیجه با اینکه در SGD در هر تکرار از بخش از داده‌ها استفاده می‌شود و

گرایان تقریبی است ولی در بلندمدت و در میانگین، مسیر به سمت مینیمم تابع هزینه است.

11. ابتدا فرمول کلی را می‌نویسیم:

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} L(\theta_t; x_i, y_i)$$

به صورتی که  $\theta_t$  پارامترهای مدل در گام  $t$ ،  $\eta$  نرخ یادگیری و  $L$  تابع هزینه است

در این صورت تحت این شرایط SGD می‌تواند به مینیمم سراسری یا محلی همگرا باشد:

الف) اگر تابع هزینه محدب باشد؛ در این حالت یک مینیمم وجود دارد. اگر شرایط زیر محیا باشد:

•  $L(\theta)$  محدب و مشتق پذیر  
• گرایان‌ها دارای واریانس محدود باشند

• نرخ یادگیری با شرایط رو به رو کاهش یابد:

$$\sum_{t=1}^{\infty} \eta_t = \infty \quad \text{و} \quad \sum_{t=1}^{\infty} \eta_t^2 < \infty$$

با این شرایط  $\lim_{t \rightarrow \infty} \theta_t = \theta^*$

ب) اگر تابع هزینه محدب نباشد؛ همگرا می‌گردد اما به مینیمم سراسری نیست اما با همان شروطی توان به مینیمم محلی رسید.



$$f_j(y; a) = a y^{a-1}, \quad y \in (0,1), \quad a > 0$$

.12

$$\Rightarrow L(a; y_1, y_2, \dots, y_N) = \prod_{i=1}^N a y_i^{a-1} = a^N \prod_{i=1}^N y_i^{a-1} \quad \text{تابع likelihood}$$

$$\Rightarrow \ln L(a) = N \ln a + (a-1) \sum_{i=1}^N \ln y_i \quad \text{لوگاریتم آن:}$$

$$\frac{d}{da} \ln L(a) = \frac{N}{a} + \sum_{i=1}^N \ln y_i = 0 \Rightarrow \frac{N}{a} = - \sum_{i=1}^N \ln y_i \quad \text{شتاب نسبت به } a:$$

$$\Rightarrow a = - \frac{N}{\sum_{i=1}^N \ln y_i} \Rightarrow \hat{a}_{ML} = - \left( \frac{1}{N} \sum_{i=1}^N \ln y_i \right)^{-1}$$