

HW04 - Artificial Intelligence

Saeed Rostami

Student ID: 98106542

$$D = \{(x_i, y_i)\}$$

Linear regression: $\beta_0 + \beta_1 x_i$

۱. (۲۵ نمره) مسأله‌ی رگرسیون خطی ساده را در نظر بگیرید. در تعریف احتمالاتی این مسأله فرض می‌کنیم رابطی زیر بین x_i و y_i وجود دارد به طوری که $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ توزیع می‌شود.

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

می‌دانیم که σ و β_0, β_1 مقادیر ثابت نامنفی هستند.

الف) اثبات کنید که تخمین بیشینه درست نمایی دو پارامتر β_0 و β_1 برابر با کمینه کردن مجموع مربعات خطا است.

$$\text{①} \quad \underset{\beta_0, \beta_1}{\operatorname{argmin}} \sum_{i=1}^N (y_i - (\beta_0 + \beta_1 x_i))^2 \triangleq \underset{\beta_0, \beta_1}{\operatorname{argmin}} \sum_{i=1}^N (y_i - \beta_0 - \beta_1 x_i)^2 \quad \text{①}$$

کمترین مربعات

$$\xrightarrow{\text{MLE}} \text{②} \quad \underset{\beta_0, \beta_1}{\operatorname{argmax}} P([y, x] | \beta_0, \beta_1) \xrightarrow{x_i \text{ are iid}} \underset{\beta_0, \beta_1}{\operatorname{argmax}} \prod_{i=1}^N P(y_i, x_i | \beta_0, \beta_1)$$

$$\xrightarrow{\text{Negative log}} \underset{\beta_0, \beta_1}{\operatorname{argmin}} - \sum_{i=1}^N \log(P(y_i, x_i | \beta_0, \beta_1))$$

$$\Rightarrow P(y_i, x_i | \beta_0, \beta_1) = P(\epsilon_i = y_i - \beta_0 - \beta_1 x_i) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2} (y_i - \beta_0 - \beta_1 x_i)^2}$$

$$\Rightarrow \underset{\beta_0, \beta_1}{\operatorname{argmin}} \sum_{i=1}^N \underbrace{\frac{1}{2\sigma^2}}_{\text{constant}} (y_i - \beta_0 - \beta_1 x_i)^2 = \sum_{i=1}^N \log\left(\frac{1}{\sqrt{2\pi}\sigma}\right) \quad \text{constant}$$

$$\triangleq \underset{\beta_0, \beta_1}{\operatorname{argmin}} \sum_{i=1}^N (y_i - \beta_0 - \beta_1 x_i)^2 \quad \text{②}$$

$$\begin{aligned} \text{①} = \text{②} &\rightarrow \text{معادلات} \\ \left\{ \begin{aligned} \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}} \end{aligned} \right. \end{aligned}$$

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad c_i \equiv (x_i - \bar{x}) \rightarrow \beta_1 = \frac{1}{S_{xx}} \sum_{i=1}^n c_i (y_i - \bar{y}) \quad (ب)$$

$$= \frac{1}{S_{xx}} \left\{ c_1 \left[Y_1 - \frac{1}{n} [Y_1 + \dots + Y_n] \right] + c_2 \left[Y_2 - \frac{1}{n} [Y_1 + \dots + Y_n] \right] + \dots + c_n \left[Y_n - \frac{1}{n} [Y_1 + \dots + Y_n] \right] \right\}$$

$$= Y_1 \left\{ \frac{1}{S_{xx}} \left[(x_1 - \bar{x}) - \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x}) \right] \right\} + \dots + Y_n \left\{ \frac{1}{S_{xx}} \left[(x_n - \bar{x}) - \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x}) \right] \right\}$$

$$= \sum_{i=1}^n Y_i \left\{ \frac{1}{S_{xx}} \left[(x_i - \bar{x}) - \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x}) \right] \right\} \Rightarrow \boxed{\hat{\beta}_1 = \sum_{i=1}^n \frac{Y_i}{S_{xx}} (x_i - \bar{x})}$$

$$\Rightarrow \hat{\beta}_1 \sim N \left(\underbrace{\sum_{i=1}^n (\beta_0 + \beta_1 x_i)}_{S_{xx}}, \sum_{i=1}^n \sigma^2 \frac{1}{S_{xx}^2} (x_i - \bar{x})^2 \right) \quad (I)$$

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (x_i^2 - 2\bar{x}x_i + \bar{x}^2) = \sum_{i=1}^n x_i^2 - 2\bar{x}n\frac{1}{n} \sum_{i=1}^n x_i + \sum_{i=1}^n \bar{x}^2$$

$$= \sum_{i=1}^n x_i^2 - n\bar{x}^2 = \sum_{i=1}^n x_i^2 - n\bar{x} \frac{1}{n} \sum_{i=1}^n x_i = \boxed{\sum_{i=1}^n x_i (x_i - \bar{x})} \quad (II)$$

$$\xrightarrow{(II)(I)} \mu(\hat{\beta}_1) = \frac{\cancel{\beta_0 \sum (x_i - \bar{x})}}{S_{xx}} + \frac{\beta_1 \sum x_i (x_i - \bar{x})}{S_{xx}} = \beta_1 \quad \text{اریب بودن}$$

$$\rightarrow \boxed{\hat{\beta}_1 \sim N \left(\beta_1, \frac{\sigma^2}{S_{xx}} \right) \triangleq N \left(\beta_1, \frac{\sigma^2}{\sum (x_i - \bar{x})^2} \right)}$$

برای اثبات توزیع $\hat{\beta}_0$ ، مراحل تقریباً مشابه با ما می‌بینید، مگر این است که در اینجا n را نزنیم.

$$\boxed{\hat{\beta}_0 \sim N \left(\beta_0, \frac{\sigma^2}{n} \left(\frac{\sum x_i^2}{\sum (x_i - \bar{x})^2} \right) \right)}$$

$$\tilde{\beta}_1 = \frac{\sum \gamma_i y_i}{\sum \gamma_i x_i} \quad \text{such that } \sum_i \gamma_i = 0 \quad (1)$$

⇒ we derive from previous part that $\hat{\beta}_1 = \sum \frac{y_i}{s_{xx}} (x_i - \bar{x})$

$$\Rightarrow \hat{\beta}_1 = \frac{1}{s_{xx}} \sum y_i (x_i - \bar{x}) \xrightarrow{s_{xx} = \sum x_i (x_i - \bar{x})} \hat{\beta}_1 = \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x}) x_i}$$

$$\gamma_i = x_i - \bar{x} \quad \text{به معنی از این خانواده است}$$

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x}) y_i}{s_{xx}} = \frac{\sum (x_i - \bar{x}) (\beta_0 + \beta_1 x_i + \epsilon_i)}{s_{xx}} \quad (2)$$

$$\Rightarrow \hat{\beta}_1 = \frac{E(\beta_1 s_{xx})}{s_{xx}} + \frac{E(\beta_0 \sum (x_i - \bar{x}))}{s_{xx}} + \frac{E(\sum (x_i - \bar{x}) \epsilon_i)}{s_{xx}} = \beta_1$$

نابرابری با صحت

$$\text{var}(\tilde{\beta}_1) = \frac{1}{s_{xx}^2} \text{var} \left[\sum \gamma_i [\beta_0 + \beta_1 x_i + \epsilon_i] \right] \quad (3)$$

$$\xrightarrow[\text{independent variable}]{\text{independent}} \text{var}(\tilde{\beta}_1) = \frac{1}{s_{xx}^2} \left[\underbrace{\text{var}(\sum \gamma_i \beta_0)}_{\alpha} + \underbrace{\text{var}(\sum \gamma_i \beta_1 x_i)}_{\beta} + \underbrace{\text{var}(\sum \gamma_i \epsilon_i)}_{s_{xx} \sigma^2} \right]$$

$$\Rightarrow \frac{1}{s_{xx}^2} [\alpha + \beta + \sigma^2 s_{xx}] \quad \text{where } \alpha, \beta > 0$$

$$\Rightarrow \text{var}(\tilde{\beta}_1) > \text{var}(\hat{\beta}_1) \quad \checkmark$$

۲. (۲۰ نمره) فرض کنید قصد داشته باشیم مساله‌ی رگرسیون چند متغیره را در نظر بگیریم. تابع هزینه‌ای که باید کمینه شود به فرم زیر خواهد بود.

$$\min_W F(W) = \lambda W^T W + \|XW - Y\|_2^2$$

الف) اگر بخواهیم این مساله را با الگوریتم Stochastic Gradient Descent حل کنیم، شبه کد آن را بنویسید.

ب) حال فرض کنید تعریف کنیم

$$W_1 = \operatorname{argmin}_W L(W)$$

$$W_\gamma = \operatorname{argmin}_W L(W) + \lambda W^T W$$

که $L(W)$ یک تابع نامنفی است. اثبات کنید که $\|W_\gamma\|_2 \leq \|W_1\|_2$ و ارتباط آن را با فرمول بندی مساله بیان کنید.

A) pseudocode * W is augmented

1) initialize $W := 0^{n+1}$

2) for iteration $t \in [1, \dots, T]$

2.1 draw random example with replacement (x_i)

2.2 compute loss by $F(W, x_i)$

2.3 compute gradient $\Delta W = -\frac{\partial F(W, x_i)}{\partial W}$

2.4 update parameter: $W := W + \gamma \Delta W$
coefficient

$$B) W_1 = \operatorname{argmin}_W L(W) \quad (I) \quad W_2 = \operatorname{argmin}_W L(W) + \lambda W^T W \quad (II)$$

proof that $\|W_2\|_2 \leq \|W_1\|_2$:

→ we prove this, by contradiction, suppose $\|W_2\|_2 \gg \|W_1\|_2$

$$\text{if } \|W_2\|_2 \gg \|W_1\|_2 \Rightarrow \lambda \|W_2\|_2^2 \gg \lambda \|W_1\|_2^2 \Rightarrow \lambda W_2^T W_2 > \lambda W_1^T W_1 \quad (1)$$

also we know from (I) that $\forall W_0 \in D(W) \quad L(W_1) \leq L(W_0)$

$$\Rightarrow L(W_1) \leq L(W_2) \quad (2)$$

using (1) (2) → $L(W_1) + \lambda W_1^T W_1 \leq L(W_2) + \lambda W_2^T W_2 \Rightarrow W_2 \neq \operatorname{argmin}_W L(W) + \lambda W^T W$

$$\checkmark \quad \|W_2\|_2 \leq \|W_1\|_2$$

← contradiction

۳. (۱۰ نمره) یک دیتاست چند متغیره را در نظر بگیرید، به این معنی که $x_i \in \mathbb{R}^p$ که $p > 1$ است و $y_i \in \mathbb{R}$ است. فرض کنید مشاهده کرده‌ایم که یکی از ضرایب محاسبه شده یک مقدار خیلی بزرگ منفی نسبت به باقی متغیرها پیدا کرده است کدام یک از گزاره‌های زیر صحیح است؟ توضیح دهید.

- این ویژگی تاثیر زیادی روی مدل دارد و باید حفظ شود.
- این ویژگی تاثیر زیادی روی مدل ندارد و باید ایگنور شود.
- نمی‌توان بدون در دست داشتن اطلاعات بیشتر در مورد این ویژگی نظر داد.

نویسنده نظر دقیقی در این رابطه دارد:

۱. ممکن است این فریب بزرگ ناشی از $overfitting$ مدل باشد، رتبه $attribute$ ما منفیه نباشد.

۲. ممکن است این $attribute$ ، یک ویژگی نه‌حس دهنه‌کننده باشد در عمل $prediction$ رتبه این ویژگی تعیین شود.

۴. (۱۵ نمره) با ارائه دلیل صحیح یا غلط بودن هر یک از گزاره‌های زیر را ثابت کنید.

- اگر $bias$ زیاد است اضافه کردن تعداد داده‌های آموزش کمک زیادی به کم کردن بایاس نمی‌کند.
- کم کردن خطای مدل روی داده‌های آموزش منجر به کاهش خطای مدل روی داده‌های تست می‌شود.
- افزایش پیچیدگی مدل رگرسیون همواره منجر به کاهش خطای مدل روی داده‌ی آموزش و افزایش خطای مدل روی داده‌ی تست می‌شود.

① غلط. افزایش تعداد داده‌های آموزش باعث می‌شود خطا، روی تعدادی از داده‌ها منبسط شود، در نتیجه $bias$ کم می‌شود. رجاسیت مدل افزایش می‌یابد.

② غلط. لزوماً افزایش رتبه نمی‌باشد. ممکن است با کم کردن خطای مدل روی داده‌ی آموزش، $overfitting$ رخ دهد که منجر به افزایش خطای روی داده‌ی تست می‌شود.

③ غلط! ممکن است در ابتدا، مدل رگرسیون بسیار ساده باشد. در این صورت با افزایش پیچیدگی، کاهش خطای روی داده‌ی تست را در پی دارد.

HEART ATTACK	EXERCISES	SMOKES	MALE	CHEST PAIN	PATIENT ID
yes	yes	no	yes	yes	۱
yes	no	yes	yes	yes	۲
yes	no	yes	no	no	۳
no	yes	no	yes	no	۴
yes	yes	yes	no	yes	۵
no	yes	yes	yes	no	۶

آلگوریتم یاد کردن درخت تصمیم گیری:

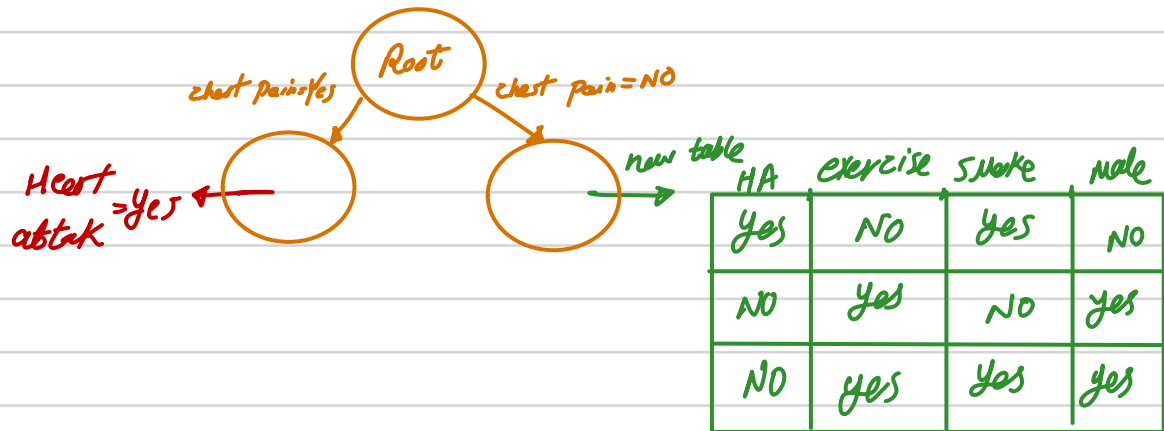
- ① مناسبی $entropy$ را
- ② مناسبی IG را برای $attribute$
- انتخاب $attribute$ با بیشترین IG

Repeat step 2 ③

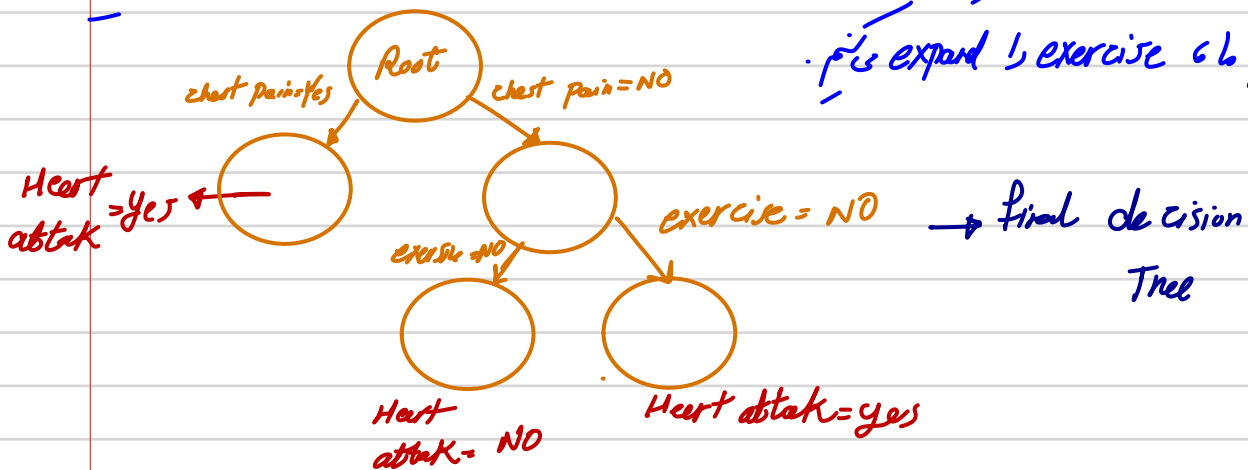
الف) با استفاده از این داده‌ها درخت تصمیم گیری پیش بینی حمله قلبی را تشکیل دهید.
ب) درخت به دست آمده را به صورت تعدادی گزاره‌ی تصمیم گیری ترجمه کنید.

$$H(\text{initial}) = - \left[\frac{4}{6} \log_2 \frac{4}{6} + \frac{2}{6} \log_2 \frac{2}{6} \right] = 0.918$$

→ برای مرحله ۱) راه مفیدی این است که $expand$ کردن کدام $attribute$ ، برآیندی کمتر را ایجاد کند. که در این مثال واضحاً با $expand$ کردن $chest\ pain$ ، کمترین برآیندی/افلا تا به بین ترین $information\ gain$ را خواهیم داشت.



که دوباره مرحله ۱) را انجام می‌دهیم. به $male$ را $expand$ کنیم، چه $exercise$ هر دو بیشترین IG خواهد داشت. در نتیجه ما، $exercise$ را $expand$ می‌کنیم.



- ① با آرنزد، $chest\ pain$ داشته باشد، قطعاً متنبه به حمله قلبی می‌شود.
- ② آرنزد $chest\ pain$ نداشته باشد و فعالیت ورزشی نداشته باشد متنبه به حمله قلبی می‌شود.
- ③ آرنزد $chest\ pain$ نداشته باشد و فعالیت ورزشی داشته باشد متنبه به حمله قلبی نمی‌شود.

۶. (۱۰ نمره) نشان دهید هر دسته بند دودویی به فرم $\{0, 1\}^d \mapsto \{0, 1\}$ می‌تواند به صورت یک درخت تصمیم‌گیری به عمق حداکثر $d + 1$ با گره‌های به فرم $(x_i = 0?)$ برای یک $i \in \{1, \dots, d\}$ پیاده‌سازی شود.

به بهترین حالت زمانی رخ می‌دهد که entropy در هر سطح، تمام attributes ما، زیر ۱ شود.
(یعنی با احتمال $\frac{1}{2}$ ، $y = 0$ ، با احتمال $\frac{1}{2}$ ، $y = 1$ ، احتمال شود). در این صورت به d گره نیاز داریم.
تا به $leaf$ node برسیم. پس یک محق دیگر که بیش از d گره نیاز به یک محق برای $prediction$ خودی رسم. به با محق d ی‌ترن $decision$ tree راحت.
حداکثر