

به نام خدا



دانشگاه صنعتی شریف
دانشکده مهندسی کامپیوتر

هوش مصنوعی

تمرین ۴: ML, Decision Tree, Regression

مهراد میلانلو

۹۹۱۰۵۷۷۵

سوال ۱

(الف)

می‌دانیم:

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

$$\epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$$\Rightarrow y_i = f(x_i; \beta) + \epsilon_i$$

$$\Rightarrow y_i = f(x_i; \beta) + \mathcal{N}(0, \sigma^2)$$

اکنون از رابطه‌ی بالا امیدریاضی می‌گیریم. خواهیم داشت:

$$\mathbb{E}[y_i | x_i] = \mathbb{E}[f(x_i; \beta) + \mathcal{N}(0, \sigma^2)] = \mathbb{E}[f(x_i; \beta)] + \underbrace{\mathbb{E}[\mathcal{N}(0, \sigma^2)]}_0 = f(x_i; \beta)$$

$$\Rightarrow p(y_i | x_i, \beta, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(y_i - f(x_i; \beta))^2\right\}$$

اکنون برای به‌دست آوردن پارامترهای β_0, β_1 مقدار likelihood را با فرض اینکه (x_i, y_i) ها *i.i.d* هستند تشکیل می‌دهیم:

$$p(\mathcal{D} | \beta, \sigma^2) = p(Y | X, \beta, \sigma^2) = \prod_{i=1}^N p(y_i | x_i, \beta, \sigma^2)$$

$$\hat{\beta} = \arg \max_{\beta} = \arg \max_{\beta} \prod_{i=1}^N p(y_i | x_i, \beta, \sigma^2)$$

برای راحتی کار، مقدار log احتمال likelihood را بیشینه می‌کنیم. چون تابع log صعودی است، جواب بهینه تغییری نمی‌کند. بنابراین خواهیم داشت:

$$\begin{aligned} \hat{\beta} &= \arg \max_{\beta} \ln p(\mathcal{D} | \beta, \sigma) = \arg \max_{\beta} \ln \prod_{i=1}^N p(y_i | x_i, \beta, \sigma^2) \\ &= \arg \max_{\beta} \sum_{i=1}^N \ln p(y_i | x_i, \beta, \sigma^2) = \arg \max_{\beta} \sum_{i=1}^N \ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(y_i - f(x_i; \beta))^2\right\} \right) \\ &= \arg \max_{\beta} \left(-N \ln \sigma - \frac{N}{2} \ln 2\pi - \underbrace{\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f(x_i; \beta))^2}_{\text{sum of square errors}} \right) \end{aligned}$$

همان‌طور که واضح است، بیشینه کردن مقدار likelihood پارامترهای β_0, β_1 معادل است با کمینه کردن مجموع مربعات خطای تخمین. بنابراین حکم موردنظر اثبات می‌شود.

(ب)

همان‌طور که در قسمت (الف) به‌دست آوردیم، برای بیشینه کردن مقدار likelihood پارامترهای β_0, β_1 کافیه مقدار

$$\sum_{i=1}^N (y_i - f(x_i; \beta))^2$$

را کمینه کنیم. به همین دلیل نسبت به β از آن مشتق می‌گیریم و برابر صفر قرار می‌دهیم:

$$Z = \sum_{i=1}^N (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\begin{cases} \frac{\partial Z}{\partial \beta_0} = -2 \sum_{i=1}^N (y_i - \beta_0 - \beta_1 x_i) = 0 \\ \frac{\partial Z}{\partial \beta_1} = -2 \sum_{i=1}^N x_i (y_i - \beta_0 - \beta_1 x_i) = 0 \end{cases} \quad (1)$$

$$\stackrel{(1)}{\Rightarrow} \sum_{i=1}^N y_i - N\beta_0 - \beta_1 \sum_{i=1}^N x_i = 0 \Rightarrow N\bar{y} - N\beta_0 - N\beta_1 \bar{x} = 0$$

$$\Rightarrow \boxed{\beta_0 = \bar{y} - \beta_1 \bar{x}}$$

$$\stackrel{(1)}{\Rightarrow} - \sum_{i=1}^N y_i x_i + \beta_0 \sum_{i=1}^N x_i + \beta_1 \sum_{i=1}^N x_i^2 = 0 \Rightarrow \sum_{i=1}^N y_i x_i - (\bar{y} - \beta_1 \bar{x}) \sum_{i=1}^N x_i - \beta_1 \sum_{i=1}^N x_i^2 = 0$$

$$\Rightarrow \underbrace{\sum_{i=1}^N y_i x_i - N\bar{y}\bar{x}}_{\sum_{i=1}^N (y_i - \bar{y})(x_i - \bar{x})} + N\beta_1 \bar{x}^2 - \beta_1 \sum_{i=1}^N x_i^2 = 0 \Rightarrow \sum_{i=1}^N (y_i - \bar{y})(x_i - \bar{x}) = \beta_1 \underbrace{\left(\sum_{i=1}^N x_i^2 - N\bar{x}^2 \right)}_{\sum_{i=1}^N (x_i - \bar{x})^2}$$

$$\Rightarrow \boxed{\beta_1 = \frac{\sum_i y_i (x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2}}$$

$$\begin{cases} \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 = \frac{\sum_i y_i (x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \end{cases} \quad (2)$$

$$\stackrel{(2)}{\Rightarrow} \mathbb{E}[\hat{\beta}_1] = \frac{1}{\sum_i (x_i - \bar{x})^2} \cdot \mathbb{E}\left[\sum_{i=1}^N y_i (x_i - \bar{x})\right] = \frac{1}{\sum_i (x_i - \bar{x})^2} \cdot \sum_{i=1}^N (x_i - \bar{x}) \mathbb{E}[y_i]$$

$$\stackrel{y_i \sim \mathcal{N}(\beta_0 + \beta_1 x_i, \sigma^2)}{y_i = \beta_0 + \beta_1 x_i + \epsilon_i} \mathbb{E}[\hat{\beta}_1] = \frac{1}{\sum_i (x_i - \bar{x})^2} \cdot \underbrace{\left(\sum_{i=1}^N \beta_0 (x_i - \bar{x}) + \sum_{i=1}^N \beta_1 x_i (x_i - \bar{x}) \right)}_{=0} = \frac{\beta_1}{\sum_i (x_i - \bar{x})^2} \cdot \sum_{i=1}^N (x_i - \bar{x}) (x_i - \bar{x})$$

$$\Rightarrow \boxed{\mathbb{E}[\hat{\beta}_1] = \beta_1}$$

$$\stackrel{(2)}{\Rightarrow} \mathbb{E}[\hat{\beta}_0] = \mathbb{E}[\bar{y} - \hat{\beta}_1 \bar{x}] = \mathbb{E}[\bar{y}] - \bar{x} \mathbb{E}[\hat{\beta}_1] = \mathbb{E}\left[\frac{\sum_i y_i}{N}\right] - \beta_1 \bar{x}$$

$$= \frac{1}{N} \mathbb{E}\left[\sum_{i=1}^N (\beta_0 + \beta_1 x_i)\right] - \beta_1 \bar{x} = \beta_0 + \beta_1 \bar{x} - \beta_1 \bar{x}$$

$$\Rightarrow \boxed{\mathbb{E}[\hat{\beta}_0] = \beta_0}$$

بنابراین مقدار امیدریاضی تخمین‌گر متغیرهای β_0, β_1 برابر با خود متغیرهاست و ثابت می‌شود تخمین‌های به‌دست آمده نااریب (Unbiased) هستند. اکنون برای به‌دست آوردن واریانس تخمین‌گرها داریم:

(دقت کنید مانند بخش الف، (x_i, y_i) ها را $i.i.d$ در نظر می‌گیریم که فرض معقولی است. همچنین x_i ها برخلاف y_i ها، متغیر Stochastic نیستند و در برخورد با آن‌ها مانند ضریب رفتار می‌کنیم. از آنجایی که این تخمین‌گرها برابر با جمع تعدادی متغیر نرمال و ضرایب هستند، واضح است که از توزیع نرمال پیروی می‌کنند.)

$$\begin{aligned}
 \xrightarrow{(\gamma)} Var(\hat{\beta}_\gamma) &= Var\left(\frac{\sum_i y_i(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^\gamma}\right) = \frac{1}{(\sum_i (x_i - \bar{x})^\gamma)^\gamma} \cdot \sum_{i=1}^N (x_i - \bar{x})^\gamma Var(y_i) \\
 &\Rightarrow \boxed{Var(\hat{\beta}_\gamma) = \frac{\sigma^\gamma}{\sum_i (x_i - \bar{x})^\gamma}} \\
 \xrightarrow{(\gamma)} \hat{\beta}_\gamma &= \bar{y} - \frac{\bar{x} \sum_i y_i(x_i - \bar{x})}{\sum_i x_i(x_i - \bar{x})} = \frac{\bar{y} \sum_i x_i^\gamma - N\bar{y}\bar{x}^\gamma - \bar{x} \sum_i x_i y_i + N\bar{y}\bar{x}^\gamma}{\sum_i x_i(x_i - \bar{x})} = \frac{\bar{y} \sum_i x_i^\gamma - \bar{x} \sum_i x_i y_i}{\sum_i x_i(x_i - \bar{x})} \\
 &= \frac{(\sum_i y_i)(\sum_i x_i^\gamma) - N\bar{x} \sum_i x_i y_i}{N \sum_i x_i(x_i - \bar{x})} \\
 \Rightarrow Var(\hat{\beta}_\gamma) &= \frac{1}{N^\gamma (\sum_i x_i(x_i - \bar{x}))^\gamma} \cdot (N(\sum_i x_i^\gamma)^\gamma - N^\gamma \bar{x}^\gamma \sum_i x_i^\gamma) \sigma^\gamma = \frac{N(\sum_i x_i^\gamma)(\sum_i x_i(x_i - \bar{x}))}{N^\gamma (\sum_i x_i(x_i - \bar{x}))^\gamma} \sigma^\gamma \\
 &\Rightarrow \boxed{Var(\hat{\beta}_\gamma) = \frac{\sigma^\gamma \sum_i x_i^\gamma}{N \sum_i (x_i - \bar{x})^\gamma}} \\
 \Rightarrow \hat{\beta}_\gamma &\sim \mathcal{N}(\beta_\gamma, \frac{\sigma^\gamma}{\sum_i (x_i - \bar{x})^\gamma}), \quad \hat{\beta}_\gamma \sim \mathcal{N}(\beta_\gamma, \frac{\sigma^\gamma \sum_i x_i^\gamma}{N \sum_i (x_i - \bar{x})^\gamma})
 \end{aligned}$$

(پ)

قرار می‌دهیم: $\gamma_i = x_i - \bar{x}$

آن‌گاه خواهیم داشت:

$$\sum_i \gamma_i = \sum_i (x_i - \bar{x}) = \sum_i x_i - \sum_i \bar{x} = \sum_i x_i - n\bar{x} = 0$$

همچنین:

$$\begin{aligned}
 \hat{\beta}_\gamma &= \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^\gamma} = \frac{\sum_i (x_i - \bar{x})y_i}{\sum_i (x_i - \bar{x})^\gamma} = \frac{\sum_i (x_i - \bar{x})y_i}{\sum_i (x_i - \bar{x})x_i} \\
 &\xrightarrow{\gamma_i = x_i - \bar{x}} \boxed{\hat{\beta}_\gamma = \frac{\gamma_i y_i}{\gamma_i x_i}}
 \end{aligned}$$

بنابراین تخمین‌گر MLE $\hat{\beta}_\gamma$ عضوی از خانواده‌ی $\tilde{\beta}_\gamma$ است.

(ت)

کافیست ثابت کنیم امیدریاضی هر تخمین‌گر عضو این خانواده، برابر با خود متغیر است. بنابراین داریم:

$$\begin{aligned}
 \mathbb{E}[\tilde{\beta}_\gamma] &= \mathbb{E}\left[\frac{\sum_i \gamma_i y_i}{\sum_i \gamma_i x_i}\right] = \frac{1}{\sum_i \gamma_i x_i} \cdot \mathbb{E}[\sum_i \gamma_i y_i] = \frac{1}{\sum_i \gamma_i x_i} \cdot \sum_i \gamma_i \mathbb{E}[y_i] \\
 &\xrightarrow{y_i \sim \mathcal{N}(\beta_\gamma + \beta_\gamma x_i, \sigma^\gamma)} \mathbb{E}[\tilde{\beta}_\gamma] = \frac{\beta_\gamma \sum_i \gamma_i + \beta_\gamma \sum_i \gamma_i x_i}{\sum_i \gamma_i x_i} \\
 &\Rightarrow \boxed{\mathbb{E}[\tilde{\beta}_\gamma] = \beta_\gamma}
 \end{aligned}$$

(ث)

ابتدا واریانس $\tilde{\beta}_1$ را محاسبه می‌کنیم:

$$\begin{aligned} \text{Var}(\tilde{\beta}_1) &= \text{Var}\left(\frac{\sum_i \gamma_i y_i}{\sum_i \gamma_i x_i}\right) = \frac{1}{\left(\sum_i \gamma_i x_i\right)^2} \cdot \text{Var}\left(\sum_i \gamma_i y_i\right) = \frac{1}{\left(\sum_i \gamma_i x_i\right)^2} \cdot \sum_i \gamma_i^2 \text{Var}(y_i) \\ &\Rightarrow \boxed{\text{Var}(\tilde{\beta}_1) = \frac{\sigma^2 \sum_i \gamma_i^2}{\left(\sum_i \gamma_i x_i\right)^2}} \end{aligned}$$

اکنون می‌خواهیم ثابت کنیم:

$$\begin{aligned} \text{Var}(\hat{\beta}_1) &\leq \text{Var}(\tilde{\beta}_1) \\ \Leftrightarrow \text{Var}(\hat{\beta}_1) &= \frac{\sigma^2}{\sum_i (x_i - \bar{x})^2} \leq \text{Var}(\tilde{\beta}_1) = \frac{\sigma^2 \sum_i \gamma_i^2}{\left(\sum_i \gamma_i x_i\right)^2} \Leftrightarrow \frac{1}{\sum_i (x_i - \bar{x})^2} \leq \frac{\sum_i \gamma_i^2}{\left(\sum_i \gamma_i x_i\right)^2} \\ &\Leftrightarrow \left(\sum_i \gamma_i x_i\right)^2 \leq \left(\sum_i \gamma_i^2\right) \left(\sum_i (x_i - \bar{x})^2\right) \end{aligned}$$

طبق نامساوی کوشی-شوارتز می‌دانیم:

$$\left(\sum_i \alpha_i \beta_i\right)^2 \leq \left(\sum_i \alpha_i^2\right) \left(\sum_i \beta_i^2\right)$$

بنابراین تنها کافیت ثابت کنیم:

$$\sum_i (\gamma_i x_i) = \sum_i \gamma_i (x_i - \bar{x})$$

برای این هم داریم:

$$\sum_i \gamma_i (x_i - \bar{x}) = \sum_i \gamma_i x_i - \bar{x} \underbrace{\sum_i \gamma_i}_{=0} = \sum_i \gamma_i x_i$$

پس حکم ثابت می‌شود زیرا:

$$\left(\sum_i \gamma_i x_i\right)^2 = \left(\sum_i \gamma_i (x_i - \bar{x})\right)^2 \leq \left(\sum_i \gamma_i^2\right) \left(\sum_i (x_i - \bar{x})^2\right)$$

خانواده‌ی $\tilde{\beta}_1$ یک خانواده‌ی تخمین‌گر خطی نااریب است. همچنین نشان دادیم تخمین‌گر MLE $\hat{\beta}_1$ عضوی از این خانواده‌است. کم‌تر بودن واریانس تخمین‌گر MLE نسبت به هر عضو دلخواه دیگری از این خانواده، نشان می‌دهد در بین تخمین‌گرهای خطی نااریب، تخمین‌گر MLE دقت بهتری دارد.

سوال ۲

(الف)

در این سوال می‌خواهیم تابع $F(W)$ را کمینه کنیم. شبه‌کد الگوریتم Stochastic Gradient Descent این‌گونه است:

```
1 stochastic_gradient_descent(F: Cost function,  $\eta$ : Learning rate, num_iter: Number of Iterations, D
   : Data):
2     Initialize:  $W \leftarrow W$ .
3     for iteration in range(num_iter):
4         S  $\leftarrow$  a random subset from Data
5         for sample in S:
6              $W = W - \eta \nabla_W F_{\text{sample}}(W)$ 
7     return W
```

Listing ۱: Stochastic Gradient Descent

دقت کنید تفاوت الگوریتم SGD با GD عادی این است که در هر مرحله از الگوریتم، به جای این که مقدار خروجی را روی تمام داده‌ها آپدیت کنیم، یک زیرمجموعه‌ی رندوم از داده‌ها را گرفته و روی آن این عمل را انجام می‌دهیم. در این جا تابع هزینه‌ی ما برابر است با:

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

بنابراین گرادین این تابع که در شبه کد بالا نیز استفاده شده است برابر است با:

$$\nabla F(W) = 2\lambda W + 2X^T(XW - Y)$$

(ب)

با استفاده از برهان خلف فرض می‌کنیم این گونه نباشد و داشته باشیم:

$$\|W_1\| < \|W_2\|$$

با توجه به این فرض، طبق تعریف واضح است که:

$$L(W_1) \leq L(W_2)$$

پس خواهیم داشت:

$$\left\{ \begin{array}{l} W_1^T W_1 < W_2^T W_2 \\ L(W_1) \leq L(W_2) \end{array} \right\} \xRightarrow{\lambda > 0} L(W_1) + \lambda W_1^T W_1 < L(W_2) + \lambda W_2^T W_2 \xRightarrow{W_2 = \arg \min L(W) + \lambda W^T W} Impossible.$$

بنابراین با رسیدن به تناقض حکم ثابت می‌شود.

در اصل عملی که انجام می‌دهیم، Regularization برای جلوگیری از مشکل Overfit است. به طوری که با اضافه کردن جمله‌ی $\lambda W^T W$ و کوچک‌تر کردن پارامترها، اثر آن‌ها را کم‌تر می‌کنیم. که این موضوع را در بالا نشان دادیم.

سوال ۳

نمی‌توان بدون در دست داشتن اطلاعات بیشتر در مورد این ویژگی نظر داد.

در واقع مقایسه‌ی بزرگی یک ضریب نسبت به ضرایب دیگر معیار مناسبی برای بررسی Overfit نیست. زیرا ممکن است یک ضریب به دلیل تاثیر و ارزش بالای feature مربوطه بزرگ باشد و یا از طرفی ممکن است برعکس آن رخ دهد. ممکن است واحد یک feature در مقایسه با feature های دیگر طوری باشد که ضریب آن بسیار بزرگ شود. بنابراین با داده‌ی محدود این سوال نمی‌توان نظر داد و اطلاعات بیشتری نیاز است.

سوال ۴

• اگر bias زیاد است اضافه کردن تعداد داده‌های آموزش کمک زیادی به کم کردن بایاس نمی‌کند. : درست

افزایش تعداد داده‌های آموزش در صورتی ممکن است به کم کردن بایاس کمک کند که واریانس مدل بالا باشد. در واقع bias به این معنی است که مدل در حالت Underfit قرار دارد. این مشکل نیز با اضافه کردن داده حل نمی‌شود و احتمالاً مدل بیش از حد ساده است و باید مدل را تغییر بدهیم.

- کم کردن خطای مدل روی داده‌های آموزش منجر به کاهش خطای مدل روی داده‌های تست می‌شود. : نادرست
ممکن است مدل ما در حالت Overfit قرار بگیرد. همان‌طور که در «سوال ۳» نیز به‌نوعی این موضوع را بررسی کردیم، ممکن است Overfit باعث بشود خطای مدل روی داده‌های آموزش بسیار کم باشد اما روی داده‌های تست عملکرد فاجعه‌باری داشته باشد. بنابراین با پیچیدگی زیاد مدل و کاهش خطای کارکرد آن روی داده‌های آموزش لزومی ندارد خطای آن روی داده‌های تست کاهش پیدا کند.
- افزایش پیچیدگی مدل رگرسیون همواره منجر به کاهش خطای مدل روی داده‌ی آموزش و افزایش خطای مدل روی داده‌ی تست می‌شود. : نادرست
اگر مدل بیش از حد ساده باشد و در حالت Underfit قرار داشته باشد، باید برای کاهش خطای آن روی داده‌های تست پیچیدگی آن را افزایش بدهیم و لزومی ندارد که خطای آن روی داده‌ی تست افزایش پیدا کند.

سوال ۵

(الف)

می‌خواهیم درخت تصمیم‌گیری پیش‌بینی حمله‌ی قلبی را بسازیم. بنابراین ریشه‌ی درخت Attack Heart است. برای سادگی نوشتن علائم در رابطه‌ها، اختصارهای زیر را استفاده می‌کنیم:

$$\begin{cases} HeartAttack \rightarrow HA \\ Exercises \rightarrow E \\ Smokes \rightarrow S \\ Male \rightarrow M \\ ChestPain \rightarrow CP \end{cases}$$

برای محاسبه‌ی آنتروپی داریم:

$$H(X) = - \sum_{i=1}^k P(X = x_i) \log_2 P(X = x_i)$$

ابتدا آنتروپی Attack Heart را محاسبه می‌کنیم:

$$H(HA) = -\frac{4}{9} \log_2\left(\frac{4}{9}\right) - \frac{5}{9} \log_2\left(\frac{5}{9}\right) \approx 0.92$$

اکنون برای یافتن node Decision بعدی، برای تمام node ها، مقدار Gain Information یا IG را حساب می‌کنیم. برای محاسبه‌ی IG داریم:

$$IG(X) = H(Y) - H(Y|X)$$

$$H(Y|X) = - \sum_{j=1}^v P(X = x_j) \sum_{i=1}^k P(Y = y_i | X = x_j) \log_2 P(Y = y_i | X = x_j)$$

بنابراین با محاسبه‌ی IG ها در این مرحله داریم:

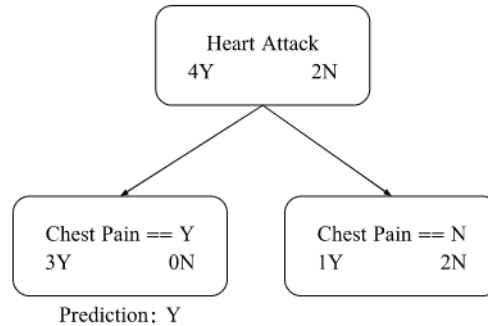
$$IG(E) = H(HA) - H(HA|E) = H(HA) + \frac{1}{3} \left(\frac{1}{3} \log_2\left(\frac{1}{3}\right) + \frac{1}{3} \log_2\left(\frac{1}{3}\right) \right) + \frac{1}{3} \left(1 \log_2(1) + 0 \log_2(0) \right) \approx 0.92 - 0.66 = 0.26$$

$$IG(S) = H(HA) - H(HA|S) = H(HA) + \frac{1}{3} \left(\frac{3}{4} \log_2\left(\frac{3}{4}\right) + \frac{1}{4} \log_2\left(\frac{1}{4}\right) \right) + \frac{1}{3} \left(\frac{1}{3} \log_2\left(\frac{1}{3}\right) + \frac{1}{3} \log_2\left(\frac{1}{3}\right) \right) \approx 0.92 - 0.87 = 0.05$$

$$IG(M) = H(HA) - H(HA|M) = H(HA) + \frac{2}{3} \left(\frac{1}{3} \log_2 \left(\frac{1}{3} \right) + \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right) + \frac{1}{3} (1 \log_2(1) + 0 \log_2(0)) \approx 0.92 - 0.66 = 0.26$$

$$IG(CP) = H(HA) - H(HA|CP) = H(HA) + \frac{1}{4} (1 \log_2(1) + 0 \log_2(0)) + \frac{1}{4} \left(\frac{1}{3} \log_2 \left(\frac{1}{3} \right) + \frac{2}{3} \log_2 \left(\frac{2}{3} \right) \right) \approx 0.92 - 0.46 = 0.46$$

از آنجایی که فیچر Pain Chest بین بقیه‌ی فیچرها، IG ماکسیمم را دارد، آن را قرار می‌دهیم. تا به این‌جا شکل Tree Decision این‌گونه است:



اکنون برای پیدا کردن فیچر بعدی، میان نمونه‌هایی با مقدار $ChestPain == No$ مجدداً IG ها را حساب می‌کنیم:

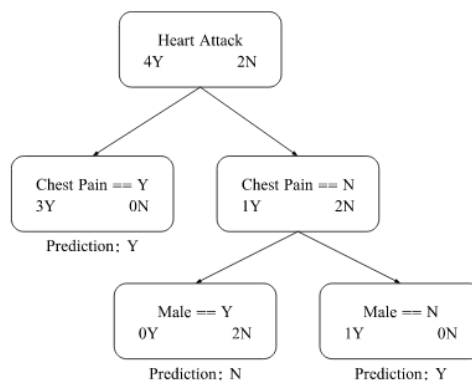
$$H(HA) = -\frac{4}{6} \log_2 \left(\frac{4}{6} \right) - \frac{2}{6} \log_2 \left(\frac{2}{6} \right) \approx 0.92$$

$$IG(E) = H(HA) - H(HA|E) = H(HA) + \frac{2}{3} (1 \log_2(1) + 0 \log_2(0)) + \frac{1}{3} (1 \log_2(1) + 0 \log_2(0)) \approx 0.92 - 0.00 = 0.92$$

$$IG(S) = H(HA) - H(HA|S) = H(HA) + \frac{2}{3} \left(\frac{1}{3} \log_2 \left(\frac{1}{3} \right) + \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right) + \frac{1}{3} (1 \log_2(1) + 0 \log_2(0)) \approx 0.92 - 0.66 = 0.26$$

$$IG(M) = H(HA) - H(HA|M) = H(HA) + \frac{2}{3} (1 \log_2(1) + 0 \log_2(0)) + \frac{1}{3} (1 \log_2(1) + 0 \log_2(0)) \approx 0.92 - 0.00 = 0.92$$

بنابراین هر دو فیچر Male, Exercises می‌توانند فیچر بعدی ما باشند. ما فیچر Male را انتخاب می‌کنیم و در نهایت درخت تصمیم به شکل زیر می‌شود:



(ب)

از ریشه به ترتیب شروع به نوشتن گزاره‌های تصمیم‌گیری می‌کنیم:

```

۱ if ChestPain == True:
۲     HeartAttack = False
۳ if ChestPain == False:
  
```



```
۴     if Exercises == True:
۵         HeartAttack = False
۶     if Exercises == False:
۷         HeartAttack = True
```

سوال ۶

طبق صورت سوال، فرض کنید فیچرها با لیبل‌های $\{x_1, \dots, x_d\}$ باشند. در بدترین حالت درختی که تشکیل می‌دهیم یک درخت دودویی کامل است که در عمق i ام تصمیم‌گیری بر اساس فیچر x_i انجام می‌شود. بنابراین عمق درخت دقیقاً $d + 1$ خواهد بود. در عمقی که node ها لیبل x_d دارند، برای بچه‌هایشان می‌توان دقیقاً مقدار خروجی را قرار داد. اکنون یک دسته‌بندی دلخواه h را در نظر بگیرید. برای هر n تایی مرتب از 0 و 1 ها مانند d به‌طوری که $d \in D_h$ باشد، واضح است که d متناظر با یک مسیر از ریشه به یکی از برگ‌های درخت ساخته‌شده است. در عمق آخر نیز با توجه به مقدار x_1, \dots, x_d ، خروجی برگ متناظر با مسیر، خروجی تابع را تعیین می‌کند.
