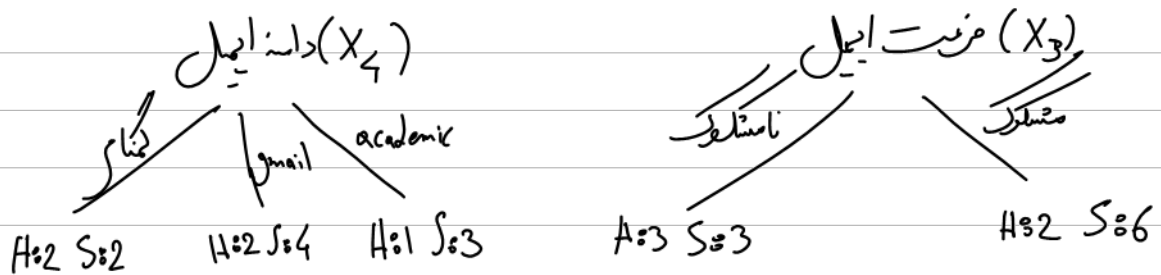
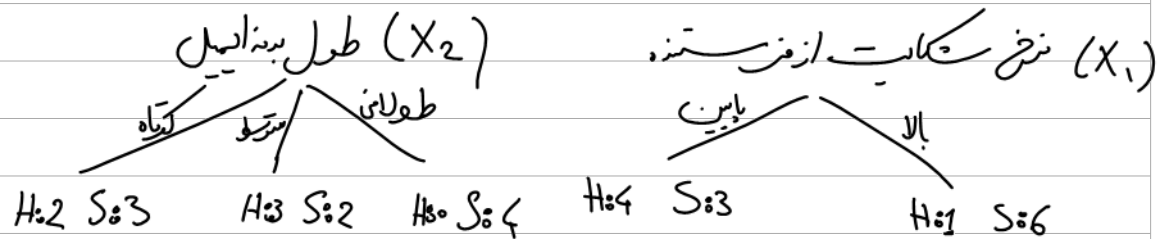


# ۱) مسأله اول

چهار ویژگی (feature) داریم:



بایستی information gain را حساب کنیم.

طبق آنکه داریم:

$$\arg \max_i I(Y, X_i) = \arg \min_i H(Y | X_i)$$

حال می‌خواهیم از  $H(Y | X_i)$  محاسبه می‌کنیم.

$$H(Y) = - \sum_y P(Y=y) \log P(Y=y)$$

$$H(Y | X_1) = \sum_x P(X_1=x) H(Y | X_1=x)$$

$$= \frac{-(\frac{3}{7} \log \frac{3}{7} + \frac{4}{7} \log \frac{4}{7}) - (\frac{1}{7} \log \frac{1}{7} + \frac{6}{7} \log \frac{6}{7})}{2} = \frac{+0.29 + 0.17}{2} = 0.23$$

$$H(Y | X_2) = \frac{3(\frac{0}{4} \log \frac{0}{4} + \frac{4}{4} \log \frac{4}{4}) - 5(\frac{3}{5} \log \frac{3}{5} + \frac{2}{5} \log \frac{2}{5}) - 1(\frac{2}{5} \log \frac{2}{5} + \frac{3}{5} \log \frac{3}{5})}{14}$$

$$\Rightarrow H(Y | X_2) = \frac{+0.29 \times 5 + 4 \times 0}{14} = 0.18$$

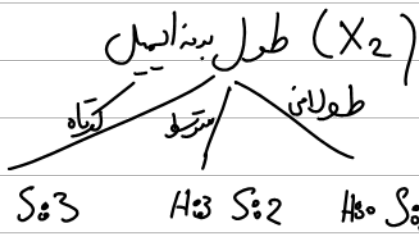
$$H(Y | X_3) = \frac{-8(\frac{2}{8} \log \frac{2}{8} + \frac{6}{8} \log \frac{6}{8}) - 6(\frac{3}{6} \log \frac{3}{6} + \frac{3}{6} \log \frac{3}{6})}{14}$$

$$H(Y | X_3) = \frac{+8 \times 0.29 + 6 \times 0.3}{14} = 0.26$$

$$H(Y|X_4) = \frac{-4(\frac{1}{4} \log \frac{1}{4} + \frac{3}{4} \log \frac{3}{4}) - 6(\frac{2}{6} \log \frac{2}{6} + \frac{4}{6} \log \frac{4}{6}) - 4(\frac{2}{4} \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4})}{14}$$

$$H(Y|X_4) = \frac{4 \times 0.24 + 6 \times 0.27 + 0.3 \times 4}{14} = 0.27$$

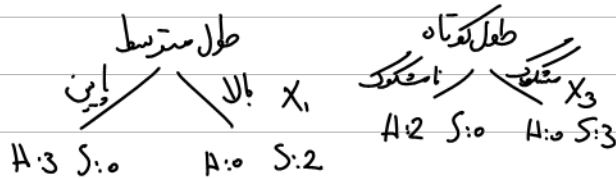
پس  $X_2$  چون از همه  $H(Y|X_i)$  کوچکتری دارد را انتخاب می‌کنیم.



حال  $X_1$  و  $X_3$  و  $X_4$  را داریم. بابر information gain محاسبه می‌کنیم.

اگر کسی دست‌کم می‌رست راست کامل می‌زد پس بالا دارد. پس نمی‌توانیم با هیچ کدام از feature ها

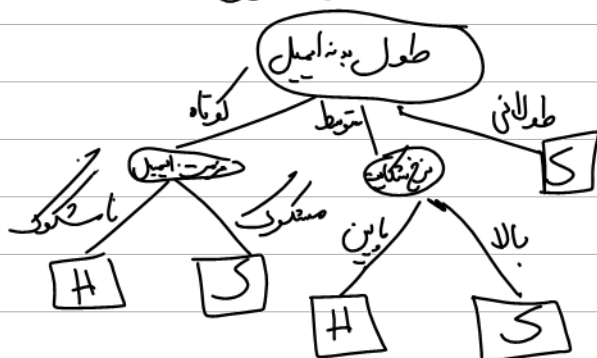
به  $I(X; Y) > 0$  برسیم. اما در صورتی که اگر دست‌کم می‌رست برای قسمت طول متوسط  $X_1$  و برای طول کوتاه



می‌بینیم که این دو feature در این دو سیر به گزینش داده‌ها را جدا می‌کنند که  $I(X; Y)$  ما کم می‌شود

در اینجا به نرم  $H:0, S:n$  یا  $H:n, S:0$  تبدیل می‌شود که می‌دانیم  $H(Y|X_i)$  در این حالت 0 است؛

که سبب max شدن information gain می‌شود. در جهت تصمیم نهایی داریم:



مست دوم

حال داده های موجود را با توجه به درخت تصمیم پیش می کنیم.

شماره	نرخ شکایت از فرستنده	طول بدنه ایمیل	فرمت ایمیل	دامنه ایمیل	تشخیص
۱	پایین	کوتاه	مشکوک	academic	نامشکوک
۲	بالا	کوتاه	نامشکوک	gmail	مشکوک
۳	بالا	کوتاه	مشکوک	گمنام	مشکوک
۴	بالا	متوسط	مشکوک	گمنام	مشکوک
۵	پایین	متوسط	مشکوک	academic	نامشکوک
۶	بالا	طولانی	نامشکوک	gmail	مشکوک

جدول ۲-۱

داده اول  $\xrightarrow{\text{طول بدنه ایمیل}} \text{کوتاه} \xrightarrow{\text{نرخ شکایت}} \text{مشکوک} \leftarrow \text{Spam}$

داده دوم  $\xrightarrow{\text{طول بدنه ایمیل}} \text{کوتاه} \xrightarrow{\text{نرخ شکایت}} \text{نامشکوک} \leftarrow \text{Ham}$

داده سوم  $\xrightarrow{\text{طول بدنه ایمیل}} \text{کوتاه} \xrightarrow{\text{نرخ شکایت}} \text{مشکوک} \leftarrow \text{Spam}$

داده چهارم  $\xrightarrow{\text{طول بدنه ایمیل}} \text{متوسط} \xrightarrow{\text{نرخ شکایت}} \text{بالا} \leftarrow \text{Spam}$

داده پنجم  $\xrightarrow{\text{طول بدنه ایمیل}} \text{متوسط} \xrightarrow{\text{نرخ شکایت}} \text{پایین} \leftarrow \text{Ham}$

داده ششم  $\xrightarrow{\text{طول بدنه ایمیل}} \text{طولانی} \leftarrow \text{Spam}$