

سوال ۱:  $\arg \max_i I(Y, X_i) = \arg \max_i [H(Y) - H(Y|X_i)] = \arg \min_i H(Y|X_i)$

$$H(Y) = - \sum_y P(Y=y) \log_2 P(Y=y), \quad H(Y|X_i) = \sum_x P(X_i=x) H(Y|X_i=x)$$

First Level:

$$\begin{aligned} H(Y| \text{نرخ شکایت از فرستنده}) &= - \frac{1}{V} \left[ \frac{3}{V} \log \frac{3}{V} + \frac{4}{V} \log \frac{4}{V} \right] - \frac{1}{V} \left[ \frac{6}{V} \log \frac{6}{V} + \frac{1}{V} \log \frac{1}{V} \right] \\ &= - \frac{1}{V} \left[ \frac{3}{V} \times 1,72 + \frac{4}{V} \times 0,8 \right] - \frac{1}{V} \left[ \frac{6}{V} \times 0,72 + \frac{1}{V} \times 2,8 \right] = \frac{1}{V} \times 0,98 + \frac{1}{V} \times 0,188 = 0,78 \end{aligned}$$

$$\begin{aligned} H(Y| \text{طول بدنه ایمیل}) &= - \frac{5}{14} \left[ \frac{3}{5} \log \frac{3}{5} + \frac{2}{5} \log \frac{2}{5} \right] - \frac{5}{14} \left[ \frac{2}{5} \log \frac{2}{5} + \frac{3}{5} \log \frac{3}{5} \right] - \frac{4}{14} \left[ \frac{4}{4} \log \frac{4}{4} + \frac{0}{4} \log \frac{0}{4} \right] \\ &= - \frac{5}{14} \left[ \frac{3}{5} \times 0,73 + \frac{2}{5} \times 1,32 \right] - \frac{5}{14} \left[ \frac{2}{5} \times 1,32 + \frac{3}{5} \times 0,73 \right] - \frac{4}{14} [1 \times 0 + 0] = 2 \times \frac{5}{14} \times 0,96 = 0,68 \end{aligned}$$

$$\begin{aligned} H(Y| \text{وقت ایمیل}) &= - \frac{8}{14} \left[ \frac{6}{8} \log \frac{6}{8} + \frac{2}{8} \log \frac{2}{8} \right] - \frac{6}{14} \left[ \frac{3}{6} \log \frac{3}{6} + \frac{3}{6} \log \frac{3}{6} \right] \\ &= - \frac{8}{14} \left[ \frac{3}{4} \times 0,41 + \frac{1}{4} \times 2 \right] - \frac{6}{14} \left[ \frac{1}{2} \times 1 + \frac{1}{2} \times 1 \right] = \frac{4}{V} \times 0,8 + \frac{3}{V} \times 1 = 0,88 \end{aligned}$$

$$\begin{aligned} H(Y| \text{دامنه ایمیل}) &= - \frac{6}{14} \left[ \frac{4}{6} \log \frac{4}{6} + \frac{2}{6} \log \frac{2}{6} \right] - \frac{4}{14} \left[ \frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4} \right] - \frac{4}{14} \left[ \frac{2}{4} \log \frac{2}{4} + \frac{2}{4} \log \frac{2}{4} \right] \\ &= - \frac{6}{14} \left[ \frac{2}{3} \times 0,58 + \frac{1}{3} \times 1,58 \right] - \frac{4}{14} \left[ \frac{3}{4} \times 0,41 + \frac{1}{4} \times 2 \right] - \frac{4}{14} \left[ \frac{1}{2} \times 1 + \frac{1}{2} \times 1 \right] \\ &= - \frac{3}{V} \times 0,91 - \frac{2}{V} \times 0,8 - \frac{2}{V} \times 1 = 0,90 \end{aligned}$$

کمترین مقدار متعلق به طول بدنه ایمیل می باشد پس آن انتخاب می شود.

Second Level: ایمیل ها با طول بدنه کوتاه.

$$\begin{aligned} H(Y| \text{نرخ شکایت از فرستنده}) &= - \frac{2}{5} \left[ \frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} \right] - \frac{3}{5} \left[ \frac{2}{3} \log \frac{2}{3} + \frac{1}{3} \log \frac{1}{3} \right] \\ &= - \frac{2}{5} \left[ \frac{1}{2} \times 1 + \frac{1}{2} \times 1 \right] - \frac{3}{5} \left[ \frac{2}{3} \times 0,58 + \frac{1}{3} \times 1,58 \right] = \frac{2}{5} + \frac{3}{5} \times 0,91 = 0,946 \end{aligned}$$

$$H(Y| \text{وقت ایمیل}) = - \frac{3}{5} \left[ \frac{3}{3} \log \frac{3}{3} + \frac{0}{3} \log \frac{0}{3} \right] - \frac{2}{5} \left[ \frac{2}{2} \log \frac{2}{2} + \frac{0}{2} \log \frac{0}{2} \right] = 0$$

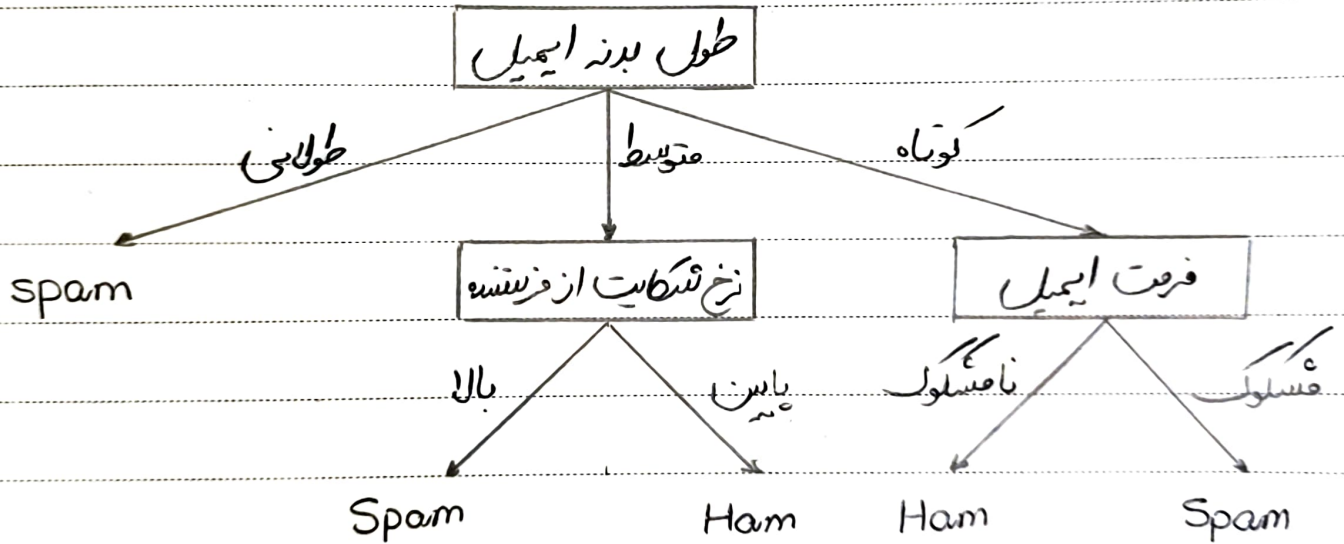
قطعاً از این کمتر نمی شود. پس وقت ایمیل را انتخاب می کنیم.

Second Level: متوسط

$H(Y) = -\frac{3}{8} \left[ \frac{2}{3} \log \frac{2}{3} + \frac{1}{3} \log \frac{1}{3} \right] - \frac{2}{8} \left[ \frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} \right] = 0.875$   
 (نرخ شکایت از فرستنده)  $H(Y)$

فقط از این کمتر نمی شود پس نرخ شکایت از فرستنده را انتخاب می کنیم.

Second Level: چوبه هم spam هستند یازاری به شاخه جدید نیست → طولانی



Spam ۴ Ham ۱ Spam ۳ Ham ۲ Spam ۱