



طول عنوان	تعداد کلمات	موضوع	تعداد خطوط	نوع فایل ضمیمه	برچسب
15	45	تبلیغات	10	Pdf	اسپم
20	30	خرید	7	Pdf	اسپم
10	25	خرید	6	Txt	غیر اسپم
30	60	فروش	8	Txt	اسپم
25	50	فروش	9	Txt	اسپم
18	40	خرید	7	Png	غیر اسپم
12	20	خرید	8	Txt	غیر اسپم
22	35	تبلیغات	9	Pdf	اسپم
17	28	فروش	7	Pdf	اسپم
14	22	خرید	6	Txt	غیر اسپم
27	55	تبلیغات	12	Png	اسپم
23	48	فروش	7	Png	اسپم
19	37	تبلیغات	8	Pdf	غیر اسپم
16	30	خرید	10	Png	غیر اسپم

➤ با کمک الگوریتم بیز ساده پیش‌بینی کنید که آیا ایمیل با مشخصات زیر اسپم خواهد بود یا خیر؟ توزیع ویژگی‌های پیوسته را نرمال در نظر بگیرید.

(18, 45, خرید, Pdf)

فاطمه شاکری، دانشکده ریاضی و علوم کامپیوتر، دانشگاه صنعتی امیرکبیر

Pdf

$$P(\text{Pdf} | \text{اسپم}) = 4/8 = 1/2$$

$$P(\text{Pdf} | \text{غیر اسپم}) = 1/6$$

خرید

$$P(\text{خرید} | \text{اسپم}) = 1/8$$

$$P(\text{خرید} | \text{غیر اسپم}) = 5/6$$

طول عنوان

$$\mu(\text{اسپم}) = (15 + 20 + 30 + 25 + 22 + 17 + 27 + 23)/8 = 22.375$$

$$\text{variance}(\text{اسپم}) =$$

$$\text{sqrt}((54.391 + 5.641 + 58.141 + 6.891 + 0.141 + 28.891 + 21.391 + 0.391)/7)$$

$$= 5.012$$

$$P(\text{اسپم} | \text{طول عنوان} = 18) = \frac{e^{-\frac{1}{2} \left(\frac{x - \mu_i}{\sigma_i} \right)^2}}{\sigma_i \sqrt{2\pi}} = 0.05438$$

$$\mu(\text{غير اسپم}) = (10+18+12+14+19+16)/6 = 14.833$$

$$\text{Variance (غير اسپم)} =$$

$$\text{sqrt}((23.357+10.029+8.025+0.694+17.363+1.361)/5) = 3.488$$

$$P(\text{غير اسپم} \mid \text{طول عنوان} = 18) = \frac{e^{-\frac{1}{2}\left(\frac{x-\mu_i}{\sigma_i}\right)^2}}{\sigma_i\sqrt{2\pi}} = 0.07573$$

تعداد کلمات

$$\mu(\text{اسپم}) = (45+30+60+50+35+28+55+48)/8 = 43.875$$

$$\text{variance(اسپم)} = \text{sqrt}((1.266 + 192.516 + 260.016 + 37.516 + 78.766 + 252.016 + 123.766 + 17.016)/7)$$

$$= 11.728$$

$$P(\text{اسپم} \mid \text{تعداد کلمات} = 45) = \frac{e^{-\frac{1}{2}\left(\frac{x-\mu_i}{\sigma_i}\right)^2}}{\sigma_i\sqrt{2\pi}} = 0.03386$$

$$\mu(\text{غير اسپم}) = (25 + 40 + 20 + 22 + 37 + 30)/6 = 29$$

$$\text{Variance (غير اسپم)} = \text{sqrt}((16 + 121 + 81 + 49 + 64 + 1)/5) = 8.149$$

$$P(\text{غير اسپم} \mid \text{تعداد کلمات} = 45) = \frac{e^{-\frac{1}{2}\left(\frac{x-\mu_i}{\sigma_i}\right)^2}}{\sigma_i\sqrt{2\pi}} = 0.00712$$

تعداد خطوط

$$\mu(\text{اسپم}) = (10+7+8+9+9+7+12+7)/8 = 8.625$$

$$\text{variance(اسپم)} = \text{sqrt}((1.891 + 2.641 + 0.391 + 0.141 + 0.141 + 2.641 + 2.641 + 11.391)/7)$$

$$= 1.768$$

$$P(\text{اسپم} \mid \text{تعداد خطوط} = 15) = \frac{e^{-\frac{1}{2}\left(\frac{x-\mu_i}{\sigma_i}\right)^2}}{\sigma_i\sqrt{2\pi}} = 0.00033$$

$$\mu(\text{غير اسپم}) = (6 + 7 + 8 + 6 + 8 + 10)/6 = 7.5$$

$$\text{Variance (غیر اسپم)} = \text{sqrt}((2.25 + 0.25 + 0.25 + 2.25 + 0.25 + 6.25)/5) = 1.517$$

$$P(\text{غیر اسپم} \mid \text{تعداد خطوط} = 15) = \frac{e^{-\frac{1}{2}\left(\frac{x-\mu_i}{\sigma_i}\right)^2}}{\sigma_i\sqrt{2\pi}} = 0.00000129489$$

P(اسپم | (18, 45 , خرید , 15, Pdf)) relative to:

$$8/14 * 0.05438 * 0.03386 * 1/8 * 0.00033 * \frac{1}{2} = 2.17011 * 10^{(-8)}$$

P(غیر اسپم | (18, 45 , خرید , 15, Pdf)) relative to:

$$6/14 * 0.07573 * 0.00712 * 5/6 * 0.00000129489 * 1/6 =$$

$$4.15596 * 10^{(-11)}$$

بنابراین تحت شرایط خواسته شده پیشبینی این است که ایمیل اسپم باشد.