

Lecture 9

Multiple word spam filter

Suppose we check for N words: w_1, w_2, \dots, w_N .

Define the 0-1 random variables $X_i = \mathbb{1}\{\text{message has word } w_i\}$

Suppose $X_1 = a_1, X_2 = a_2, \dots, X_N = a_N$, where a_i are either 0 or 1.

Assume each word appears in a message independent of the other words

$$P(X_1 = a_1, X_2 = a_2, \dots, X_N = a_N \mid \text{spam})$$

$$= P(X_1 = a_1 \mid \text{spam})P(X_2 = a_2 \mid \text{spam}) \cdots P(X_N = a_N \mid \text{spam})$$

independence (Naive Bayes)

Example: Consider the earlier example:

	<i>spam</i>	<i>ham</i>
	1500	3672
<i>meeting</i>	16	153
<i>pharmacy</i>	621	0
<i>money</i>	125	31
<i>Digipen</i>	0	1892

$$P(\text{spam}) = \frac{1500}{1500 + 3672} = 0.29$$

$$P(\text{ham}) = 0.71$$

Use smoothing, with smoothing parameters $(\alpha, \beta) = (1, 2)$.

(a) Email message has the words w_1, w_3 , and w_4 , but not the word w_2

$$P(\text{spam} \mid X_1 = 1, X_2 = 0, X_3 = 1, X_4 = 1) = ?$$

(b) Email message has the words w_1, w_3 , but not the words w_2 or w_4

$$P(\text{spam} \mid X_1 = 1, X_2 = 0, X_3 = 1, X_4 = 0) = ?$$

$$P(s | X_1 = 1, X_2 = 0, X_3 = 1, X_4 = 1) = \frac{P(X_1 = 1, X_2 = 0, X_3 = 1, X_4 = 1 | s) P(s)}{P(X_1 = 1, X_2 = 0, X_3 = 1, X_4 = 1)}$$

$$\begin{aligned} P(X_1 = 1, X_2 = 0, X_3 = 1, X_4 = 1 | \text{spam}) &= P(X_1 = 1 | \text{spam}) P(X_2 = 0 | \text{spam}) \\ &\quad \times P(X_3 = 1 | \text{spam}) P(X_4 = 1 | \text{spam}) \\ &= \frac{16+1}{1500+2} \cdot \left(1 - \frac{62+1}{1502}\right) \cdot \frac{125+1}{1502} \cdot \frac{1}{1502} = \textcircled{1} \end{aligned}$$

$$P(X_1 = 1, X_2 = 0, X_3 = 1, X_4 = 1 | \text{ham}) = \frac{154}{3674} \times \left(1 - \frac{1}{3674}\right) \times \frac{32}{3674} \times \frac{1893}{3674} = \textcircled{2}$$

$$\begin{aligned} P(\text{spam} | X_1 = 1, X_2 = 0, X_3 = 1, X_4 = 1) &= \frac{\textcircled{1} (0.29)}{\textcircled{1} (0.29) + \textcircled{2} (0.71)} \\ &= 0.0008 \end{aligned}$$

→ label "ham"

$$P(s | X_1 = 1, X_2 = 0, X_3 = 1, X_4 = 0) = \frac{P(X_1 = 1, X_2 = 0, X_3 = 1, X_4 = 0 | s) P(s)}{P(X_1 = 1, X_2 = 0, X_3 = 1, X_4 = 0)}$$

$$\begin{aligned} &\frac{\frac{17}{1502} \cdot \left(1 - \frac{622}{1502}\right) \cdot \frac{126}{1502} \cdot \left(1 - \frac{1}{1502}\right) (0.29)}{=} \\ &= \frac{\text{(above)} + \frac{154}{3674} \left(1 - \frac{1}{3674}\right) \cdot \frac{32}{3674} \cdot \left(1 - \frac{1893}{3674}\right) (0.71)}{=} \\ &= 0.56 \end{aligned}$$

→ label "spam"

[depends on threshold]

Testing the model

We can use the following metrics to test the model:

$$\text{accuracy} = \frac{\text{correct predictions}}{\text{total}}$$

$$\text{precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}} = \frac{\text{spam predict spam}}{\text{predict spam}}$$

$$\text{recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}} = \frac{\text{spam predict spam}}{\text{total spam}}$$

Example: Consider the email predictions:

	spam	ham
predict spam	101	33
predict ham	38	704

The three evaluation metrics we can use give:

$$(a) \text{ accuracy} = \frac{101 + 704}{101 + 33 + 38 + 704} = .9189$$

$$(b) \text{ precision} = \frac{101}{101 + 33} = .7537$$

$$(c) \text{ recall} = \frac{101}{101 + 38} = .7266$$

Compact formulation of model

Goal: *pre-compute* parameters for the model, based on training data.

Suppose our spam filter keeps track of N different words. (w_1, w_2, \dots, w_N)

When a new email message arrives, it is encoded by a vector

$$\vec{a} = [a_1, a_2, \dots, a_N]$$

$a_i = 1$ if w_i is in message
 $a_i = 0$ if w_i is not

with 1's for the words that appear in the message and 0's for those that do not appear.

Notation: We will use the following facts and notation in our derivations:

✓ $\exp\{x\} = e^x$, for an easier way to display the expressions,

✓ $\sum_{k=1}^n c_k = c_1 + c_2 + \dots + c_n,$

✓ $\prod_{k=1}^n c_k = c_1 \times c_2 \times \dots \times c_n,$

✓ exponentials and logs are inverses of each other: $x = e^{\log(x)}$, $\log(e^x) = x$

✓ property of logs: $\log(a \cdot b) = \log(a) + \log(b),$

✓ property of logs: $\log(a^b) = b \log(a).$

- for a large N set of words, we need to multiply many small probabilities, leading to underflow problems.

$$\begin{aligned}
P(X_1 = a_1, \dots, X_N = a_N | \text{spam}) &= P(X_1 = a_1 | \text{spam}) \times \dots \times P(X_N = a_N | \text{spam}) \\
&= \prod_{k=1}^N P(X_k = a_k | \text{spam}) \\
&= \exp \{ \log(P(X_1 = a_1 | \text{spam}) \times \dots \times P(X_N = a_N | \text{spam})) \} \\
&= \exp \{ \log(P(X_1 = a_1 | \text{spam})) + \dots + \log(P(X_N = a_N | \text{spam})) \} \\
&= \exp \left\{ \sum_{k=1}^N \log(P(X_k = a_k | \text{spam})) \right\},
\end{aligned}$$

For a more compact way to write these probabilities, we let

$$p_{ks} = P(X_k = 1 | \text{spam}), \quad p_{kh} = P(X_k = 1 | \text{ham})$$

be the probabilities that w_k appears as a spam message / ham message.

$$\begin{aligned}
P(X_k = a_k | \text{spam}) &= P(X_k = 1 | \text{spam})^{a_k} [1 - P(X_k = 1 | \text{spam})]^{1-a_k} \\
&= p_{ks}^{a_k} (1 - p_{ks})^{1-a_k}
\end{aligned}$$

$$\begin{aligned}
&\text{if } a_k = 0 \Rightarrow P(X_k = 0 | \text{spam}) = 1 - p_{ks} = (1 - p_{ks})^{1-a_k} \\
&\text{if } a_k = 1 \Rightarrow P(X_k = 1 | \text{spam}) = p_{ks} = p_{ks}^{a_k}
\end{aligned}$$

$$\begin{aligned}
P(X_k = a_k | \text{ham}) &= P(X_k = 1 | \text{ham})^{a_k} [1 - P(X_k = 1 | \text{ham})]^{1-a_k} \\
&= p_{kh}^{a_k} (1 - p_{kh})^{1-a_k}
\end{aligned}$$

$$\log a - \log b = \log \frac{a}{b}$$

Combining into the logarithm notation:

$$\begin{aligned} & P(X_1 = a_1, \dots, X_N = a_N \mid \text{spam}) \\ &= \exp \left\{ \sum_{k=1}^N \log [p_{ks}^{a_k} (1 - p_{ks})^{1-a_k}] \right\} \\ &= \exp \left\{ \sum_{k=1}^N [a_k \log(p_{ks}) + (1 - a_k) \log(1 - p_{ks})] \right\} \\ &= \exp \left\{ \sum_{k=1}^N \left[a_k \log \left(\frac{p_{ks}}{1 - p_{ks}} \right) \right] + \sum_{k=1}^N \log(1 - p_{ks}) \right\}. \end{aligned}$$

depends on training data.

Let $y_0 = \sum_{k=1}^N \log(1 - p_{ks})$ and $y_k = \log \left(\frac{p_{ks}}{1 - p_{ks}} \right)$.

Set $\vec{X} = [X_1, \dots, X_N]$, $\vec{a} = [a_1, \dots, a_N]$ and $\vec{y} = [y_1, \dots, y_N]$:

$$P(\vec{X} = \vec{a} \mid \text{spam}) = \exp\{\vec{a} \cdot \vec{y} + y_0\}.$$

Let $z_0 = \sum_{k=1}^N \log(1 - p_{kh})$, $\vec{z} = [z_1, z_2, \dots, z_N]$ if $z_k = \log \left(\frac{p_{kh}}{1 - p_{kh}} \right)$.

$$P(\vec{X} = \vec{a} \mid \text{ham}) = \exp\{\vec{a} \cdot \vec{z} + z_0\}.$$

With y_0 , \vec{y} , z_0 , and \vec{z} *pre-computed*, we classify messages:

$$P(\text{spam} \mid \vec{X} = \vec{a}) = \frac{\exp\{\vec{y} \cdot \vec{a} + y_0\} P(\text{spam})}{\exp\{\vec{y} \cdot \vec{a} + y_0\} P(\text{spam}) + \exp\{\vec{z} \cdot \vec{a} + z_0\} P(\text{ham})}$$

Example: using the 4 words below, with smoothing $\alpha = 1$ and $\beta = 2$

	spam	ham
	1500	3672
<i>meeting</i>	16	153
<i>pharmacy</i>	621	0
<i>money</i>	125	31
<i>DigiPen</i>	0	1892

$$y_0 = \sum \log(1 - p_{ks})$$

$$y_k = \log\left(\frac{p_{ks}}{1 - p_{ks}}\right)$$

Recall $P(\text{spam}) = .29$ and $P(\text{ham}) = .71$. Then:

$$\vec{y} = \left[\log\left(\frac{\frac{17}{1502}}{1 - \frac{17}{1502}}\right), \log\left(\frac{\frac{622}{1502}}{1 - \frac{622}{1502}}\right), \log\left(\frac{\frac{126}{1502}}{1 - \frac{126}{1502}}\right), \log\left(\frac{\frac{1}{1502}}{1 - \frac{1}{1502}}\right) \right]$$

$$= [-4.47, -0.35, -2.39, -7.31]$$

$$y_0 = \log\left(1 - \frac{17}{1502}\right) + \log\left(1 - \frac{622}{1502}\right) + \log\left(1 - \frac{126}{1502}\right) + \log\left(1 - \frac{1}{1502}\right)$$

$$= -0.63$$

$$\vec{z} = \left[\log\left(\frac{\frac{154}{3674}}{1 - \frac{154}{3674}}\right), \log\left(\frac{\frac{1}{3674}}{1 - \frac{1}{3674}}\right), \log\left(\frac{\frac{32}{3674}}{1 - \frac{32}{3674}}\right), \log\left(\frac{\frac{1893}{3674}}{1 - \frac{1893}{3674}}\right) \right]$$

$$= [-3.13, -8.21, -4.73, -0.06]$$

$$z_0 = \log\left(1 - \frac{154}{3674}\right) + \log\left(1 - \frac{1}{3674}\right) + \log\left(1 - \frac{32}{3674}\right) + \log\left(1 - \frac{1893}{3674}\right)$$

$$= -0.776$$

$$\begin{aligned}
& P(\vec{X} = [1, 0, 1, 1] \mid \text{spam}) \\
&= \exp \{ [-4.47, -0.35, -2.39, -7.31] \cdot [1, 0, 1, 1] + (-0.63) \} \\
&= e^{-14.8} \quad \begin{array}{ccc} \vec{y} & \vec{a} & y_0 \end{array} \\
&= 3.7 \times 10^{-7}
\end{aligned}$$

$$\begin{aligned}
& P(\vec{X} = [1, 0, 1, 1] \mid \text{ham}) \\
&= \exp \{ [-3.13, -8.21, -4.73, -0.06] \cdot [1, 0, 1, 1] + (-0.776) \} \\
&= e^{-8.696} \quad \begin{array}{ccc} \vec{z} & \vec{a} & z_0 \end{array} \\
&= 1.67 \times 10^{-4}
\end{aligned}$$

$$\begin{aligned}
& P(\text{spam} \mid \vec{X} = [1, 0, 1, 1]) \\
&= \frac{P(\vec{X} = [1, 0, 1, 1] \mid \text{spam}) P(\text{spam})}{P(\vec{X} = [1, 0, 1, 1] \mid \text{spam}) P(\text{spam}) + P(\vec{X} = [1, 0, 1, 1] \mid \text{ham}) P(\text{ham})} \\
&= \frac{(3.7 \times 10^{-7})(.29)}{(3.7 \times 10^{-7})(.29) + (1.67 \times 10^{-4})(.71)} \\
&= 0.00090 \quad \text{close to earlier computation.}
\end{aligned}$$

$$\begin{aligned}
& P(\text{spam} \mid \vec{X} = [1, 0, 1, 0]) \\
&= \frac{\exp\{\vec{y} \cdot [1, 0, 1, 0] + y_0\}(.29)}{\exp\{\vec{y} \cdot [1, 0, 1, 0] + y_0\}(.29) + \exp\{\vec{z} \cdot [1, 0, 1, 0] + z_0\}(.71)} \\
&= .5623 \quad \text{close.}
\end{aligned}$$