



INTELLIGENT NETWORKING & SYSTEM LAB. © 2021

## **Discrete feature values**

- Consider building an email spam filter using machine learning
- We wish to classify messages according to whether they are unsolicited commercial (spam) email, or non-spam email → Automatically filter out
- Classifying emails is one example of a broader set of problems called classification
- Training set
  - ✓ A set of emails labeled as spam or non-spam
  - ✓ We will begin our construction of our spam filter by specifying the features  $x_i$  used to represent an email

## Training set

- We will represent an email via a feature vector whose length is equal to the number of words in the dictionary
- If an email contains the *i*-th word of the dictionary, then we will set x<sub>i</sub> = 1; otherwise, we let x<sub>i</sub> = 0
  - ✓ For instance, the vector is used to represent an email that contains the words 'a' and 'buy', but not 'aardvark', 'aardwolf', or 'zygmurgy'

# Training set

- Having chosen our feature vector, we now want to build a generative model p(x | y)
- To model *p*(*x* | *y*), we will make a very strong assumption
   ✓ We will assume that the *x<sub>i</sub>*'s are conditionally independent given *y*
  - $\checkmark$  This assumption is called the Naïve Bayes assumption
  - ✓ The resulting algorithm is called the Naïve Bayes classifier
- If y = 1 means spam email; 'buy' is word 2087 and 'price' is word 39831;
  - $\checkmark$  Then we are assuming that if I tell you y = 1 (that a particular piece of email is spam),
  - ✓ Then knowledge of  $x_{2087}$  (knowledge of whether 'buy' appears in the message) will have no effect on your beliefs about the value of  $x_{39831}$  (whether 'price' appears)

$$p(x_{2087} \mid y) = p(x_{2087} \mid y, x_{39831})$$

 $\checkmark$  We are only assuming that  $x_{2087}$  and  $x_{39831}$  are conditionally independent given y

## **Training set**

$$p(x_1, ..., x_{50000} | y)$$
  
=  $p(x_1 | y) p(x_2 | y, x_1) p(x_3 | y, x_1, x_2) \cdots p(x_{50000} | y, x_1, ..., x_{49999})$   
=  $p(x_1 | y) p(x_2 | y) p(x_3 | y) \cdots p(x_{50000} | y)$   
=  $\prod_{i=1}^{n} p(x_i | y)$ 

- The first equality simply follows from the usual properties of probabilities
- The second equality used the Naïve Bayes assumption
- Even though the Naïve Bayes assumption is an extremely strong assumptions, the resulting algorithm works well on many problems

Our model is parameterized by

$$\phi_{i|y=1} = p(x_i = 1 | y = 1), \ \phi_{i|y=0} = p(x_i = 1 | y = 0) \text{ and } \phi_y = p(y = 1)$$

- Given a training set  $\{(x^{(i)}, y^{(i)}); i = 1, \dots, m\}$
- Write down the joint likelihood of the data

$$L(\phi_{y}, \phi_{j|y=0}, \phi_{j|y=1}) = \prod_{i=1}^{m} p(x^{(i)}, y^{(i)})$$

• Maximizing this with respect to  $\phi_y, \phi_{j|y=0}, \phi_{j|y=1}$  gives the maximum likelihood estimates

• The maximum likelihood estimates:

$$\phi_{j|y=1} = \frac{\sum_{i=1}^{m} 1\{x_{j}^{(i)} = 1 \land y^{(i)} = 1\}}{\sum_{i=1}^{m} 1\{y^{(i)} = 1\}}$$

$$\phi_{j|y=0} = \frac{\sum_{i=1}^{m} 1\{x_{j}^{(i)} = 1 \land y^{(i)} = 0\}}{\sum_{i=1}^{m} 1\{y^{(i)} = 0\}}$$

$$\phi_{y} = \frac{\sum_{i=1}^{m} 1\{y^{(i)} = 1\}}{m}$$

In the equations above, the 'A ' symbol means 'and'

 Having fit all these parameters, to make a prediction on a new example with features x, we then simply calculate

$$p(y=1|x) = \frac{p(x|y=1)p(y=1)}{p(x)}$$
$$= \frac{\prod_{j=1}^{n} p(x_j|y=1)p(y=1)}{\prod_{j=1}^{n} p(x_j|y=1)p(y=1) + \prod_{j=1}^{n} p(x_j|y=0)p(y=0)}$$

The Naïve Bayes classifier finds the most possible state with the largest probability as the posterior probability for y

$$v = \arg \max_{y} p(y=1|X)$$
  
=  $\arg \max_{y} p(X|y=1) p(y=1)$   
=  $\arg \max_{y} \prod_{j=1}^{n} p(x_{j}|y=1) p(y=1)$ 

#### Laplace smoothing

- The Naïve Bayes algorithm as we have described it will work fairly well for many problems, but there is a simple change that makes it work much better
- In case of work "nips"
  - ✓ We had not previously seen any emails containing the word "nips"
  - ✓ "nips" did not ever appear in your training set of spam/non-spam emails
  - $\checkmark$  Assuming that "nips" was the 35000<sup>th</sup> word in the dictionary

$$\phi_{35000|y=1} = \frac{\sum_{i=1}^{m} 1\left\{x_{35000}^{(i)} = 1 \land y^{(i)} = 1\right\}}{\sum_{i=1}^{m} 1\left\{y^{(i)} = 1\right\}} = 0$$
  
$$\phi_{35000|y=0} = \frac{\sum_{i=1}^{m} 1\left\{x_{35000}^{(i)} = 1 \land y^{(i)} = 0\right\}}{\sum_{i=1}^{m} 1\left\{y^{(i)} = 0\right\}} = 0$$

#### Laplace smoothing

- Because it has never seen "nips" before in either spam or non-spam training examples, it thinks the probability of seeing it in either type of email is zero
- Hence, when trying to decide if one of these messages containing "nips" is spam

$$p(y=1|x) = \frac{\prod_{j=1}^{n} p(x_j | y=1) p(y=1)}{\prod_{j=1}^{n} p(x_j | y=1) p(y=1) + \prod_{j=1}^{n} p(x_j | y=0) p(y=0)}$$
  
= 0

Our algorithm obtains 0/0, and doesn't know how to make a prediction

# Laplace smoothing

• To avoid this problem, we can use *Laplace smoothing* 

$$\phi_{j|y=1} = \frac{\sum_{i=1}^{m} 1\{x_{j}^{(i)} = 1 \land y^{(i)} = 1\} + 1}{\sum_{i=1}^{m} 1\{y^{(i)} = 1\} + k}$$
$$\phi_{j|y=0} = \frac{\sum_{i=1}^{m} 1\{x_{j}^{(i)} = 1 \land y^{(i)} = 0\} + 1}{\sum_{i=1}^{m} 1\{y^{(i)} = 0\} + k}$$

## Example

Input features: color, type, origin

No.	Color	Туре	Origin	Stolen?
1	RED	Sports	Domestic	YES
2	RED	Sports	Domestic	NO
3	RED	Sports	Domestic	YES
4	YELLOW	Sports	Domestic	NO
5	YELLOW	Sports	Imported	YES
6	YELLOW	SUV	Imported	NO
7	YELLOW	SUV	Imported	YES
8	YELLOW	SUV	Domestic	NO
9	RED	SUV	Imported	NO
10	RED	Sports	Imported	YES

 We want to classify a Red Domestic SUV. Note there is no example of a Red Domestic SUV in our data set

INTELLIGENT NETWORKING & SYSTEM LAB. © 2021

#### Example

• Find the probabilities:

$$p(x_{j} | y = 1) = \frac{\sum_{i=1}^{m} 1\{x_{j}^{(i)} = 1 \land y^{(i)} = 1\} + 1}{\sum_{i=1}^{m} 1\{y^{(i)} = 1\} + k} \quad p(x_{j} | y = 0) = \frac{\sum_{i=1}^{m} 1\{x_{j}^{(i)} = 1 \land y^{(i)} = 0\} + 1}{\sum_{i=1}^{m} 1\{y^{(i)} = 0\} + k}$$

$$p(y = 1) = \frac{\sum_{i=1}^{m} 1\{y^{(i)} = 1\} + 1}{m + k} \quad p(y = 0) = \frac{\sum_{i=1}^{m} 1\{y^{(i)} = 0\} + 1}{m + k}$$

 $p(RED | Yes) = 0.57, \quad p(RED | No) = 0.43$  $p(SUV | Yes) = 0.29, \quad p(SUV | No) = 0.57$  $p(Domestic | Yes) = 0.43, \quad p(Domestic | No) = 0.57$ 

#### Example

Since 0.070 > 0.036, the example gets classified as "No"

 $v_1 = p(Yes) \cdot p(RED \mid Yes) \cdot p(SUV \mid Yes) \cdot p(Domestic \mid Yes) = 0.036$  $v_2 = p(No) \cdot p(RED \mid No) \cdot p(SUV \mid No) \cdot p(Domestic \mid No) = 0.070$ 





INTELLIGENT NETWORKING & SYSTEM LAB. © 2021