# Math 170E Problem Set 2

Allan Zhang

April 16, 2025

## 1  Bayes' Rule

### 1.1  Application of Bayes' rule to medical testing

Suppose there is a new test for a disease, which has a probability of 0.99 of detecting the disease if it is present in the patient (a true positive), but also has a probability of 0.03 of detecting the disease in a person who does not have it (a false positive). Additionally, suppose that the probability that a random person in the population has the disease is 0.05.

(a) What is the probability that a random person who is tested tests positive?

Let $+$ and $-$ indicate positive and negative tests. Let $H$ and $S$ represent people who do not have the disease and people who do (happy and sad lol)

Then we know that
$$P(+|S) = 0.99 \quad P(+|H) = 0.03 \quad P(S) = 0.05$$

From the law of total probability, we know that

$$P(+) = P(+|S)P(S) + P(+|H)P(H)$$

$$P(+) = 0.99 \cdot 0.05 + 0.03 \cdot 0.95$$

$$P(+) = 0.078$$

(b) What is the probability that a random person has the disease, given they test positive?

In other words, what is $P(S|+)$

Bayes' theorem states

$$P(S|+) = \frac{P(+|S)}{P(+)}P(S)$$

$$P(S|+) = \frac{0.99}{0.078} \cdot 0.05$$

$$P(S|+) = 0.635$$

## 1.2 The Naive Bayes test classifier

Suppose you would like to classify texts based on whether they are comedy or news. As training data, you take some collections of news headlines and jokes, and list all the words in all the jokes in a list C and all the words in all the news headlines in a list $N$.

Each of the lists have 1000 total words. Further, suppose that

- The word 'man' appears 50 times in $C$ and 50 times in $N$

- The word 'eats' appears 25 times in $C$ and 10 times in $N$

- The word 'pet' appears 25 times in $C$ and 10 times in $N$

- The word 'gerbil' appears 25 times in $C$ and 2 times in $N$

Now you randomly choose one of the two lists $C$ or $N$ with equal probability, then sample (independently, with replacement, in order) four words from that list, and obtain the sample:

$$x = (x1, x2, x3, x4) = (\text{'man'}, \text{'eats'}, \text{'pet'}, \text{'gerbil'})$$

(c) What is the probability $P(C|x)$ that the sample is from $C$, the words from jokes?

We know that $P(C)$ and $P(N)$ are both 0.5. Then

$$P(C|x) = \frac{P(x|C)}{P(x)} P(C)$$

We also know that $P(x) = P(x|C)P(C) + P(x|N)P(N)$

$$P(x|C) = P(\text{'man'}|C)P(\text{'eats'}|C)P(\text{'pet'}|C)P(\text{'gerbil'}|C)$$

$$P(x|C) = \frac{50}{1000} \cdot \frac{25}{1000} \cdot \frac{25}{1000} \cdot \frac{25}{1000} = 7.84 \times 10^{-7}$$

$$P(x|N) = P(\text{'man'}|N)P(\text{'eats'}|N)P(\text{'pet'}|N)P(\text{'gerbil'}|N)$$

$$P(x|N) = \frac{50}{1000} \cdot \frac{10}{1000} \cdot \frac{10}{1000} \cdot \frac{2}{1000} = 1 \times 10^{-8}$$

$$P(x) = 0.5(7.84 \times 10^{-7}) + 0.5(1 \times 10^{-8}) = 3.97 \times 10^{-7}$$

$$P(C|x) = \frac{7.84 \times 10^{-7}}{3.97 \times 10^{-7}} \cdot 0.5$$

$$P(C|x) = 0.987$$

(d) What is the probability $P(N|x)$

We know that $P(N|x) = 1 - P(C|x)$

$$P(N|x) = 1 - 0.987 = 0.013$$

(e) If you later saw this list of words in a new sentence, how would you (naively) decide whether it is a joke or news using the principle of maximum likelihood estimation?

temp