# COMP 472: Artificial Intelligence
# Machine Learning
# Naive Bayes Classification
# Application to Spam Filtering

- Russell & Norvig: Sections 12.2 to 12.6

# Today

1. Introduction to ML
2. Naïve Bayes Classification
   a. Application to Spam Filtering  **YOU ARE HERE!**
3. Decision Trees
4. ( Evaluation
5. Unsupervised Learning )
6. Neural Networks
   a. Perceptrons
   b. Multi Layered Neural Networks

# Recall

$$H_{NB} = \underset{H_i}{\arg\max} \; \frac{P(H_i) \times P(E \mid H_i)}{P(E)} = \underset{H_i}{\arg\max} \; P(H_i) \times P(E \mid H_i) = \underset{H_i}{\arg\max} \; P(H_i) \times P(<a_1, a_2, a_3, ..., a_n> \mid H_i) = \underset{H_i}{\arg\max} \; P(H_i) \times \prod_{j=1}^{n} P(a_j \mid H_i)$$

$$H_{NB} = \underset{H_i}{\arg\max} \; P(H_i) \times \prod_{j=1}^{n} P(a_j \mid H_i)$$

# Application of Naive Bayes Classification: Spam Filtering

- Task: classify e-mails (documents) into a pre-defined class
  - ex: spam / ham
  - ex: sports, recreation, politics, war, economy,...
- Given

  spam/ham,          sports/recreation/politics
  - training set of documents already classified into the correct category

**SPAM**          **HAM**

# e-mail Representation

- each e-mail is represented by a vector of feature/value:
  - feature = actual words in the e-mail
  - value = number of times that word appears in the e-mail

1000    10000

<airplane=0, banana=1, cat=5, duck=4, ..., zoo=0, class=SPAM>

<airplane=2, banana=0, cat=0, duck=8, ..., zoo=3, class=SPAM>

...

<airplane=1, banana=1, cat=5, duck=8, ..., zoo=3, class=HAM>

<airplane=1, banana=3, cat=5, duck=0, ..., zoo=6, class=HAM>

multinomial
naive bayes classifier

Strictly speaking, what this is called a Multinomial Naïve Bayes classifier, because we use the frequency of words, as opposed to just using binary values for the presence/absence of words.

naive bayes classifier    10                    word

# Naïve Bayes Algorithm

```
// 1. training
for all classes cᵢ   // ex. ham or spam                    conditional probability
    for all words wⱼ in the vocabulary
        compute
```

$$P(w_j \mid c_i) = \frac{count(w_j, c_i)}{\sum_j count(w_j, c_i)}$$

ci                8                                                                          H

spam

```
for all classes cᵢ
    compute
```

$$P(c_i) = \frac{count(documents\ in\ c_i)}{count(all\ documents)}$$

prior probability

training set          P(ham)        P(spam)

```
// 2. testing a new document D
for all classes cᵢ // ex. ham or spam
    score(cᵢ) = P(cᵢ)              priority
    for all words wⱼ in the D
        score(cᵢ) = score(cᵢ) x P(wⱼ | cᵢ)
                    hypothesis              ci hypothesis
choose c* = with the greatest score(cᵢ)
```

|  | $w_1$ | $w_2$ | $w_3$ | $w_4$ | $w_5$ | $w_6$ |
|---|---|---|---|---|---|---|
| c1 : SPAM | $p(w_1\|c_1)$ | $p(w_2\|c_1)$ | $p(w_3\|c_1)$ | $p(w_4\|c_1)$ | $p(w_5\|c_1)$ | $p(w_6\|c_1)$ |
| c2 : HAM | $p(w_1\|c_2)$ | $p(w_2\|c_2)$ | $p(w_3\|c_2)$ | $p(w_4\|c_2)$ | $p(w_5\|c_2)$ | $p(w_6\|c_2)$ |

# Example 1

- **Dataset**
  - c1: SPAM
    - doc1: "cheap meds for sale"
    - doc2: "click here for the best meds"
    - doc3: "book your trip"

  - c2: HAM
    - doc4: "cheap book sale, not meds"
    - doc5: "here is the book for you"

**SPAM**

**HAM**

- **Question:**
  - doc6: "the cheap book"
  - should it be classified as HAM or SPAM?

training set

# Example 1

**Assume**
vocabulary = {best, book, cheap, sale, trip, meds}
If not in vocabulary, ignore word

1. Training:

conditionality

- $P(best|SPAM) = 1/7$     $P(best|HAM) = 0/5$
- $P(book|SPAM) = 1/7$     $P(book|HAM) = 2/5$
- $P(cheap|SPAM) = 1/7$     $P(cheap|HAM) = 1/5$
- $P(sale|SPAM) = 1/7$     $P(sale|HAM) = 1/5$
- $P(trip|SPAM) = 1/7$     $P(trip|HAM) = 0/5$
- $P(meds|SPAM) = 2/7$     $P(meds|HAM) = 1/5$

prior

- $P(SPAM) = 3/5$     $P(HAM) = 2/5$

prior
class,  bayes

conditionality

2. Testing: "the cheap book"

- $Score(HAM) = P(HAM) \times P(cheap|HAM) \times P(book|HAM)$    HAM    $\times$ cheap    $\times$ book
- $Score(SPAM) = P(SPAM) \times P(cheap|SPAM) \times P(book|SPAM)$

# Be Careful: Smooth Probabilities

- normally: $P(w_i \mid c_j) = \dfrac{(\text{frequency of } w_i \text{ in } c_j)}{\text{total number of words in } c_j}$

- what if we have a $P(w_i \mid c_j) = 0...$?
  - ex. the word "dumbo" never appeared in the class SPAM?
  - then P("dumbo"| SPAM) = 0                0

- so if a text contains the word "dumbo", the class SPAM is completely ruled out !

- to solve this: we assume that every word always appears at least once (or a smaller value)   additive smoothing
  - ex: add-1 smoothing:                      add1,        1                     0.5

$$P(w_i \mid c_j) = \dfrac{(\text{frequency of } w_i \text{ in } c_j) + 1}{\text{total number of words in } c_j + \text{size of vocabulary}}$$

fake 1,        size of vocabulary(        )

# Smoothing add-1 smoothing

- **Assume:**
  - vocabulary V = {ball, heat, kitchen, referee, stove, the, ... }
  - |V| = 100      size of vocabulary

- **Training set:**

orginal data set

| c1: COOKING | c2: SPORTS |
|---|---|
| $doc_1$:  ... stove... kitchen... the... heat <br> $doc_2$:  ... kitchen... pasta... stove... <br> ... <br> $doc_{100000}$:  ... stove...heat... ball... | $doc_1$:  ... ball... heat... <br> $doc_2$:  ... the... referee... player... <br> ... <br> $doc_{75000}$:  goal... injury ... |

:ball:1,heat:1.............100

100extra word

original+      100extra word

# Be Careful: Use Logs

- if we really do the product of probabilities…

  - $\text{argmax}_{c_j}\ P(c_j) \prod P(w_i|c_j)$      0.xx
  - we soon have numerical underflow…
  - ex: 0.01 × 0.02 × 0.05 × …

Log       log    the ranking of hypothesis

- so instead, we add the log of the probs

  - $\text{argmax}_{c_j}\ \log(P(c_j)) + \sum \log(P(w_i|c))$    -3    -4
  - ex: log(0.01) + log(0.02) + log(0.05) + …

log    base     log 2 log 3, log 4

# Example 2

- Training set:

| c1: COOKING | c2: SPORTS |
|---|---|
| $doc_1$: ... stove... kitchen... the... heat<br>$doc_2$: ... kitchen... pasta... stove...<br>...<br>$doc_{100000}$: ... stove...heat... ball... | $doc_1$: ... ball... heat...<br>$doc_2$: ... the... referee... player...<br>...<br>$doc_{75000}$: goal... injury ... |

- Assume:
  - vocabulary V = {ball, heat, kitchen, referee, stove, the, ... }
  - |V| = 100
  - 500,000 words in Cooking
  - 300,000 words in Sports
  - 100,000 docs in Cooking
  - 75,000 docs in Sports

# Example 2

- Training – Unsmoothed / Smoothed probs:

  - $P(\text{ball}|\text{COOKING}) = \frac{10,000}{500,000}$  $\frac{??}{??}$ $\frac{10000+1}{5000000+100}$  $P(\text{ball}|\text{SPORTS}) = \frac{10,000}{300,000}$  $\frac{??}{??}$

  - $P(\text{heat}|\text{COOKING}) = \frac{255}{500,000}$  $\frac{??}{??}$  $P(\text{heat}|\text{SPORTS}) = \frac{1,800}{300,000}$  $\frac{??}{??}$

  - $P(\text{kitchen}|\text{COOKING}) = \frac{2,600}{500,000}$  $\frac{??}{??}$  log  $P(\text{kitchen}|\text{SPORTS}) = \frac{0}{300,000}$  $\frac{??}{??}$

  - $P(\text{referee}|\text{COOKING}) = \frac{0}{500,000}$  $\frac{??}{??}$  $P(\text{referee}|\text{SPORTS}) = \frac{1,50}{30\,,000}$  $\frac{??}{??}$

  - $P(\text{stove}|\text{COOKING}) = \frac{3,600}{500,000}$  $\frac{??}{??}$  $P(\text{stove}|\text{SPORTS}) = \frac{4}{300,000}$  $\frac{??}{??}$

  - $P(\text{the}|\text{COOKING}) = \frac{400,000}{500,000}$  $\frac{??}{??}$  $P(\text{the}|\text{SPORTS}) = \frac{19,000}{300,000}$  $\frac{??}{??}$

  - …

  - $P(\text{COOKING}) = \frac{100,000}{175,000}$        $P(\text{SPORTS}) = \frac{75,000}{175,000}$

- Testing: "the referee hit the ~~blue bird~~"

  - Score(COOKING)= $\log(\frac{100,000}{175,000})$ + log(P(the|COOKING)) + log(P(referee|COOKING)) + log(P(hit|COOKING)) + log(P(the|COOKING))

  - Score(SPORTS)= $\log(\frac{75,000}{175,000})$ + log(P(the|SPORTS)) + log(P(referee|SPORTS)) + log(P(hit|SPORTS)) + log(P(the| SPORTS))

# Today

1. Introduction to ML ✓
2. Naïve Bayes Classification ✓
   a. Application to Spam Filtering ✓
3. Decision Trees
4. ( Evaluation
5. Unsupervised Learning )
6. Neural Networks
   a. Perceptrons
   b. Multi Layered Neural Networks

# Up Next

1. Introduction to ML
2. Naïve Bayes Classification
    a. Application to Spam Filtering
3. **Decision Trees**
4. ( Evaluation
5. Unsupervised Learning )
6. Neural Networks
    a. Perceptrons
    b. Multi Layered Neural Networks