

CSCI 6350/7350 Notes on text classification using Naïve Bayes

1. Build a Naïve Bayes classifier to determine whether a user would be interested or not interested in a document. $\alpha=1$.

Assume the user has indicated he is interested in the following documents :

Doc 1: victory garden

Doc 2: square foot organic gardening

Doc 3: the victory garden companion

Doc 4: victory gardens for organic foods

And he is not interested in the

Doc 1: four season harvest

Doc 2: gardening bible

Doc 3: victory gardening art poster

Doc 4: victory gardening recipe

Show whether the classifier will determine the following document to be of interest or not of interest:

Grow organic food with victory garden method

2. Given the documents below that belong to 3 different categories: Artificial Intelligence(AI), Parallel Processing(PP), and Software Engineering(SE),

(1) construct a naïve Bayes classifier. $\alpha=1$.

AI:

Doc 1: new reinforcement learning method

Doc 2: machine learning is one important AI method

PP:

Doc 1: load balancing is important

Doc 2: balancing the documents on systems

SE:

Doc 1: System verification method

(2) Classify the following two documents into the correct category using the naïve Bayes classifier built.

Test doc: heuristic search based reinforcement learning method

TRAINMULTINOMIALNB(\mathbf{C}, \mathbf{D})

```
1   $V \leftarrow \text{EXTRACTVOCABULARY}(\mathbf{D})$   
2   $N \leftarrow \text{COUNTDOCS}(\mathbf{D})$   
3  for each  $c \in \mathbf{C}$   
4  do  $N_c \leftarrow \text{COUNTDOCSINCLASS}(\mathbf{D}, c)$   
5      $\text{prior}[c] \leftarrow N_c / N$   
6      $\text{text}_c \leftarrow \text{CONCATENATETEXTOFALLDOCSINCLASS}(\mathbf{D}, c)$   
7     for each  $t \in V$   
8     do  $T_{ct} \leftarrow \text{COUNTTOKENSOFTERM}(\text{text}_c, t)$   
9     for each  $t \in V$   
10    do  $\text{condprob}[t][c] \leftarrow \frac{T_{ct}+1}{\sum_{c'} (T_{ct'}+1)}$   
11 return  $V, \text{prior}, \text{condprob}$ 
```

APPLYMULTINOMIALNB($\mathbf{C}, V, \text{prior}, \text{condprob}, d$)

```
1   $W \leftarrow \text{EXTRACTTOKENSFROMDOC}(V, d)$   
2  for each  $c \in \mathbf{C}$   
3  do  $\text{score}[c] \leftarrow \log \text{prior}[c]$   
4     for each  $t \in W$   
5     do  $\text{score}[c] += \log \text{condprob}[t][c]$   
6 return  $\arg \max_{c \in \mathbf{C}} \text{score}[c]$ 
```

TRAINBERNOULLNB(\mathbb{C}, \mathbb{D})

```
1  $V \leftarrow \text{EXTRACTVOCABULARY}(\mathbb{D})$   
2  $N \leftarrow \text{COUNTDOCS}(\mathbb{D})$   
3 for each  $c \in \mathbb{C}$   
4 do  $N_c \leftarrow \text{COUNTDOCSINCLASS}(\mathbb{D}, c)$   
5    $\text{prior}[c] \leftarrow N_c / N$   
6   for each  $t \in V$   
7   do  $N_{ct} \leftarrow \text{COUNTDOCSINCLASSCONTAININGTERM}(\mathbb{D}, c, t)$   
8      $\text{condprob}[t][c] \leftarrow (N_{ct} + 1) / (N_c + 2)$   
9 return  $V, \text{prior}, \text{condprob}$ 
```

APPLYBERNOULLNB($\mathbb{C}, V, \text{prior}, \text{condprob}, d$)

```
1  $V_d \leftarrow \text{EXTRACTTERMSFROMDOC}(V, d)$   
2 for each  $c \in \mathbb{C}$   
3 do  $\text{score}[c] \leftarrow \log \text{prior}[c]$   
4   for each  $t \in V$   
5   do if  $t \in V_d$   
6     then  $\text{score}[c] += \log \text{condprob}[t][c]$   
7     else  $\text{score}[c] += \log(1 - \text{condprob}[t][c])$   
8 return  $\arg \max_{c \in \mathbb{C}} \text{score}[c]$ 
```