## Slide 1

# EECS 498-004: Introduction to Natural Language Processing

Instructor: Prof. Lu Wang

Computer Science and Engineering

University of Michigan

https://web.eecs.umich.edu/~wangluxy/

1

## Slide 2

# Project proposal (due Feb 15)

- In general, we want to see that you have a clear goal in the project. The technical details can be described in a rough manner, but in principle, you need to show what problem you want to study.
  - **Introduction**: the problem has to be well-defined. What are the input and output. Why this is an important problem to study.
  - **Related work**: put your work in context. Describe what has been done in previous work on the same or related subject. And why what you propose to do here is novel and different.
  - **Datasets**: what data do you want to use? What is the size of it? What information is contained? Why is it suitable for your task?
  - **Methodology**: what models do you want to use? You may change the model as the project goes, but you may want to indicate some type of models that might be suitable for your problem. Is it a supervised learning problem or unsupervised? What classifiers can you start with? Are you making improvements? You don't have to be crystal clear on this section, but it can be used to indicate the direction that your project goes to.
  - **Evaluation**: what metrics do you want to use for evaluating your models?
- Length: 1 page (or more if necessary). Single space if MS word is used. Or you can choose latex templates, e.g. https://www.acm.org/publications/proceedings-template.
- Grading: based on each section described above, 20 points per section. But as you can tell, they're related to each other.
- Each group member will make separate submission with all group members' names indicated.

2

## Slide 3

# Project discussions!

- See Piazza @14 to sign up for meeting times
- Happy to discuss your ideas or just brainstorm together!
- Feb 1, 10:30am-12pm
- Feb 3, 12pm-1pm
- Feb 4, 9pm-10pm

3

## Slide 4

# Outline

- Text Categorization/Classification
- Naïve Bayes
- Evaluation

4

## Slide 5

# Positive or negative movie review?

- unbelievably disappointing
- Full of zany characters and richly applied satire, and some great plot twists
- this is the greatest screwball comedy ever filmed
- It was pathetic. The worst part about it was the boxing scenes.

5

## Slide 6

# Male or female author?

1. By 1925 present-day Vietnam was divided into three parts under French colonial rule. The southern region embracing Saigon and the Mekong delta was the colony of Cochin-China; the central area with its imperial capital at Hue was the protectorate of Annam…

2. Clara never failed to be astonished by the extraordinary felicity of her own name. She found it hard to trust herself to the mercy of fate, which had managed over the years to convert her greatest shame into one of her greatest assets…

S. Argamon, M. Koppel, J. Fine, A. R. Shimoni, 2003. "Gender, Genre, and Writing Style in Formal Written Texts," Text, volume 23, number 3, pp. 321–346

6

## Text Classification

- Assigning subject categories, topics, or genres
- Spam detection
- Authorship identification
- Age/gender identification
- Language Identification
- Sentiment analysis
- …

7

7

## Text Classification: definition

- *Input*:
  - a document $d$
  - a fixed set of classes  $C = \{c_1, c_2, ..., c_J\}$

- *Output*: a predicted class $c \in C$

8

8

## Classification Methods:
## Hand-coded rules

- Rules based on combinations of words or other features
  - spam: black-list-address OR ("dollars" AND "have been selected")
- Accuracy can be high
  - If rules carefully refined by expert
- But building and maintaining these rules is expensive

9

9

## Classification Methods:
## Supervised Machine Learning

- *Input:*
  - a document $d$
  - a fixed set of classes  $C = \{c_1, c_2, ..., c_J\}$
  - A training set of $m$ hand-labeled documents $(d_1, y_1), ...., (d_m, y_m)$, $y_i$ is in C
- *Output:*
  - a learned classifier $\gamma : d \rightarrow c$

10

10

## Classification Methods:
## Supervised Machine Learning

- Any kind of classifier
  - Naïve Bayes
  - Logistic regression
  - Support-vector machines
  - k-Nearest Neighbors
  - Neural networks
  - …

11

11

## Outline

- Text Categorization/Classification
- Naïve Bayes
- Evaluation

12

12

## Naïve Bayes Classifier

13

13

## Naïve Bayes Intuition

- Simple ("naïve") classification method based on Bayes rule
- Relies on very simple representation of document
  - Bag of words

14

14

## The Bag of Words Representation

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!

| it | 6 |
| I | 5 |
| the | 4 |
| to | 3 |
| and | 3 |
| seen | 2 |
| yet | 1 |
| would | 1 |
| whimsical | 1 |
| times | 1 |
| sweet | 1 |
| satirical | 1 |
| adventure | 1 |
| genre | 1 |
| fairy | 1 |
| humor | 1 |
| have | 1 |
| great | 1 |
| ... | ... |

15

15

## The bag of words representation

$$\gamma( \quad \begin{array}{|l|l|} \hline \text{seen} & 2 \\ \hline \text{sweet} & 1 \\ \hline \text{whimsical} & 1 \\ \hline \text{recommend} & 1 \\ \hline \text{happy} & 1 \\ \hline \text{...} & \text{...} \\ \hline \end{array} \quad )=c$$

16

16

## Bayes' Rule Applied to Documents and Classes

- For a document *d* and a class *c*

$$P(c \mid d) = \frac{P(d \mid c)P(c)}{P(d)}$$

17

17

## Naïve Bayes Classifier (I)

$$c_{MAP} = \operatorname*{argmax}_{c \in C} P(c \mid d)$$

MAP is "maximum a posteriori" = most likely class

$$= \operatorname*{argmax}_{c \in C} \frac{P(d \mid c)P(c)}{P(d)}$$

Bayes Rule

$$= \operatorname*{argmax}_{c \in C} P(d \mid c)P(c)$$

Dropping the denominator

18

18

3

## Naïve Bayes Classifier (I)

$$c_{MAP} = \underset{c \in C}{\operatorname{argmax}} P(c \mid d)$$

MAP is "maximum a posteriori" = most likely class

$$= \underset{c \in C}{\operatorname{argmax}} \frac{P(d \mid c)P(c)}{P(d)}$$

Bayes Rule

$$= \underset{c \in C}{\operatorname{argmax}} P(d \mid c)P(c)$$

Dropping the denominator

Why we can do this?

19

## Naïve Bayes Classifier (II)

$$c_{MAP} = \underset{c \in C}{\operatorname{argmax}} P(d \mid c)P(c)$$

$$= \underset{c \in C}{\operatorname{argmax}} P(x_1, x_2, \ldots, x_n \mid c)P(c)$$

Document d represented as features x1..xn

20

20

## Naïve Bayes Classifier (IV)

$$c_{MAP} = \underset{c \in C}{\operatorname{argmax}} P(x_1, x_2, \ldots, x_n \mid c)P(c)$$

$O(|X|^n \bullet |C|)$ parameters

$|X|$ represents the maximum number of possible values for $x_i$

21

21

$$P(x_1, x_2, \ldots, x_n \mid c)$$

- **Bag of Words assumption**: Assume position doesn't matter
- **Conditional Independence**: Assume the feature probabilities $P(x_i \mid c_j)$ are independent given the class $c$.

$$P(x_1, \ldots, x_n \mid c) = P(x_1 \mid c) \bullet P(x_2 \mid c) \bullet P(x_3 \mid c) \bullet \ldots \bullet P(x_n \mid c)$$

22

22

## Naïve Bayes Classifier

$$c_{MAP} = \underset{c \in C}{\operatorname{argmax}} P(x_1, x_2, \ldots, x_n \mid c)P(c)$$

$$c_{NB} = \underset{c \in C}{\operatorname{argmax}} P(c) \prod_{x \in X} P(x \mid c)$$

23

23

## Applying Naive Bayes Classifiers to Text Classification

positions ← all word positions in test document

$$c_{NB} = \underset{c_j \in C}{\operatorname{argmax}} P(c_j) \prod_{i \in positions} P(x_i \mid c_j)$$

24

24

4

## Learning for Naïve Bayes Model

25

## Learning the Naïve Bayes Model

- First attempt: maximum likelihood estimates
  - simply use the frequencies in the data

$$\hat{P}(c_j) = \frac{doccount(C = c_j)}{N_{doc}}$$

$$\hat{P}(w_i \mid c_j) = \frac{count(w_i, c_j)}{\sum_{w \in V} count(w, c_j)}$$

26

## Parameter estimation

$$\hat{P}(w_i \mid c_j) = \frac{count(w_i, c_j)}{\sum_{w \in V} count(w, c_j)}$$

fraction of times word $w_i$ appears among all words in documents of class $c_j$

27

## Problem with Maximum Likelihood

- What if we have seen no training documents with the word **fantastic** and classified in the topic **positive** (**thumbs-up)**?

$$\hat{P}("fantastic" \mid positive) = \frac{count("fantastic", positive)}{\sum_{w \in V} count(w, positive)} = 0$$

- Zero probabilities cannot be conditioned away, no matter the other evidence!

$$c_{MAP} = \text{argmax}_c \, \hat{P}(c) \prod_i \hat{P}(x_i \mid c)$$

28

## Laplace (add-1) smoothing for Naïve Bayes

$$\hat{P}(w_i \mid c) = \frac{count(w_i, c) + 1}{\sum_{w \in V} \left( count(w, c) + 1 \right)}$$

$$= \frac{count(w_i, c) + 1}{\left( \sum_{w \in V} count(w, c) \right) + |V|}$$

29

## Naïve Bayes: Learning

- From training corpus, extract *Vocabulary*

- Calculate $P(c_j)$ terms
  - For each $c_j$ in $C$ do
    $docs_j \leftarrow$ all docs with class $= c_j$

$$P(c_j) \leftarrow \frac{|docs_j|}{|\text{total \# documents}|}$$

- Calculate $P(w_k \mid c_j)$ terms
  - $Text_j \leftarrow$ single doc containing all $docs_j$
  - For each word $w_k$ in *Vocabulary*
    $n_k \leftarrow$ # of occurrences of $w_k$ in $Text_j$

$$P(w_k \mid c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha \mid Vocabulary \mid}$$

30

## Naïve Bayes: Learning

- From training corpus, extract *Vocabulary*

- Calculate $P(c_j)$ terms
  - For each $c_j$ in $C$ do
    $docs_j \leftarrow$ all docs with class $= c_j$

$$P(c_j) \leftarrow \frac{|docs_j|}{|total\ \#\ documents|}$$

- Calculate $P(w_k | c_j)$ terms
  - $Text_j \leftarrow$ single doc containing all $docs_j$
  - For each word $w_k$ in *Vocabulary*
    $n_k \leftarrow$ # of occurrences of $w_k$ in $Text_j$

$$P(w_k | c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha\ |Vocabulary|}$$

A more general form: add-$\alpha$ smoothing!

31

31

## Naïve Bayes and Language Modeling

- Naïve bayes classifiers can use any sort of feature
  - URL, email address, dictionaries, network features
- But if, as in the previous slides
  - We use **only** word features
  - we use **all** of the words in the text (not a subset)
- Then
  - Naïve bayes has an important similarity to language modeling.

32

32

## Each class = a unigram language model

- Assigning each word: P(word | c)
- Assigning each sentence: P(sentence|c)=Π P(word|c)

Class *pos*

| 0.1 | I |
| 0.1 | love |
| 0.01 | this |
| 0.05 | fun |
| 0.1 | film |

| I | love | this | fun | film |
|---|------|------|-----|------|
| 0.1 | 0.1 | 0.01 | 0.05 | 0.1 |

P(sentence | pos) = 0.0000005

33

33

## Naïve Bayes as a Language Model

- Which class assigns the higher probability to s?

| Model pos | Model neg |
|-----------|-----------|
| 0.1 I | 0.2 I |
| 0.1 love | 0.001 love |
| 0.01 this | 0.01 this |
| 0.05 fun | 0.005 fun |
| 0.1 film | 0.1 film |

| | I | love | this | fun | film |
|--|---|------|------|-----|------|
| | 0.1 | 0.1 | 0.01 | 0.05 | 0.1 |
| | 0.2 | 0.001 | 0.01 | 0.005 | 0.1 |

P(s|pos) > P(s|neg)

34

34

## An Example

35

35

$$\hat{P}(c) = \frac{N_c}{N}$$

$$\hat{P}(w|c) = \frac{count(w,c)+1}{count(c)+|V|}$$

| | Doc | Words | Class |
|--|-----|-------|-------|
| Training | 1 | Chinese Beijing Chinese | c |
| | 2 | Chinese Chinese Shanghai | c |
| | 3 | Chinese Macao | c |
| | 4 | Tokyo Japan Chinese | j |
| Test | 5 | Chinese Chinese Chinese Tokyo Japan | ? |

**Priors:**
$P(c) = \frac{3}{4}$
$P(j) = \frac{1}{4}$

**Conditional Probabilities:**
P(Chinese|c) = (5+1) / (8+6) = 6/14 = 3/7
P(Tokyo|c) = (0+1) / (8+6) = 1/14
P(Japan|c) = (0+1) / (8+6) = 1/14
P(Chinese|j) = (1+1) / (3+6) = 2/9
P(Tokyo|j) = (1+1) / (3+6) = 2/9
P(Japan|j) = (1+1) / (3+6) = 2/9

**Choosing a class:**
P(c|d5) $\propto$ 3/4 * (3/7)³ * 1/14 * 1/14
≈ 0.0003

P(j|d5) $\propto$ 1/4 * (2/9)³ * 2/9 * 2/9
≈ 0.0001

36

36

## Summary: Naive Bayes is Not So Naive

- Very Fast, low storage requirements

- Robust to Irrelevant Features
    - Irrelevant Features cancel each other without affecting results
- Very good in domains with many equally important features

- Optimal if the independence assumptions hold:
    - If assumed independence is correct, then it is the Bayes Optimal Classifier for problem

- A good dependable baseline for text classification

37

37

## Outline

- Text Categorization/Classification
- Naïve Bayes
➡ - Evaluation

38

38

## Evaluation

39

39

## The 2-by-2 contingency table (or confusion matrix)

| | correct | not correct |
|---|---|---|
| selected | tp (true positive) | fp (false positive) |
| not selected | fn (false negative) | tn (true negative) |

For example,
- Which set of documents are related to the topic of NLP?
- Which set of documents are written by Shakespeare?

40

40

## The 2-by-2 contingency table

| | correct | not correct |
|---|---|---|
| selected | tp | fp |
| not selected | fn | tn |

41

41

## Precision and recall

- **Precision**: % of selected items that are correct, tp/(tp+fp)
  **Recall**: % of correct items that are selected, tp/(tp+fn)

| | correct | not correct |
|---|---|---|
| selected | tp | fp |
| not selected | fn | tn |

42

42

7

## A combined measure: F-measure or F-score

- A combined measure that assesses the P/R tradeoff is F measure (weighted harmonic mean):

$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha)\frac{1}{R}}$$

- People usually use balanced F1 measure
  - i.e., $\alpha = \frac{1}{2}$, $F = 2PR/(P+R)$

43

---

## Text Classification Evaluation

44

---

## More Than Two Classes: Sets of binary classifiers

- Dealing with any-of or multivalue classification
  - A document can belong to 0, 1, or >1 classes.

- For each class $c \in C$
  - Build a classifier $\gamma_c$ to distinguish c from all other classes $c' \in C$
- Given test doc d,
  - Evaluate it for membership in each class using each $\gamma_c$
  - d belongs to any class for which $\gamma_c$ returns true

45

---

## More Than Two Classes: Sets of binary classifiers

- One-of or multinomial classification
  - Classes are mutually exclusive: each document in exactly one class

- For each class $c \in C$
  - Build a classifier $\gamma_c$ to distinguish c from all other classes $c' \in C$
- Given test doc d,
  - Evaluate it for membership in each class using each $\gamma_c$
  - d belongs to the one class with maximum score

46

---

## Confusion matrix c

- For each pair of classes $<c_1, c_2>$ how many documents from $c_1$ were incorrectly assigned to $c_2$?
  - $c_{3,2}$: 90 wheat documents incorrectly assigned to poultry

| Docs in test set | Assigned UK | Assigned poultry | Assigned wheat | Assigned coffee | Assigned interest | Assigned trade |
|---|---|---|---|---|---|---|
| True UK | 95 | 1 | 13 | 0 | 1 | 0 |
| True poultry | 0 | 1 | 0 | 0 | 0 | 0 |
| True wheat | 10 | 90 | 0 | 1 | 0 | 0 |
| True coffee | 0 | 0 | 0 | 34 | 3 | 7 |
| True interest | 0 | 1 | 2 | 13 | 26 | 5 |
| True trade | 0 | 0 | 2 | 14 | 5 | 10 |

47

---

## Per class evaluation measures

**Recall**:
Fraction of docs in class $i$ classified correctly:

$$\frac{c_{ii}}{\sum_j c_{ij}}$$

**Precision**:
Fraction of docs assigned class $i$ that are actually about class $i$:

$$\frac{c_{ii}}{\sum_j c_{ji}}$$

**Accuracy**: (1 - error rate)
Fraction of docs classified correctly:

$$\frac{\sum_i c_{ii}}{\sum_j \sum_i c_{ij}}$$

48

## Micro- vs. Macro-Averaging

- If we have more than one class, how do we combine multiple performance measures into one quantity?
- **Macroaveraging**: Compute performance for each class, then average.
- **Microaveraging**: Collect decisions for all classes, compute contingency table, evaluate.

49

---

## Micro- vs. Macro-Averaging: Example

**Class 1**

| | Truth: yes | Truth: no |
|---|---|---|
| Classifier: yes | 10 | 10 |
| Classifier: no | 10 | 970 |

**Class 2**

| | Truth: yes | Truth: no |
|---|---|---|
| Classifier: yes | 90 | 10 |
| Classifier: no | 10 | 890 |

**Micro Avg. Table**

| | Truth: yes | Truth: no |
|---|---|---|
| Classifier: yes | 100 | 20 |
| Classifier: no | 20 | 1860 |

50

---

## Micro- vs. Macro-Averaging: Example

**Class 1**

| | Truth: yes | Truth: no |
|---|---|---|
| Classifier: yes | 10 | 10 |
| Classifier: no | 10 | 970 |

**Class 2**

| | Truth: yes | Truth: no |
|---|---|---|
| Classifier: yes | 90 | 10 |
| Classifier: no | 10 | 890 |

**Micro Avg. Table**

| | Truth: yes | Truth: no |
|---|---|---|
| Classifier: yes | 100 | 20 |
| Classifier: no | 20 | 1860 |

- Macroaveraged precision: (0.5 + 0.9)/2 = 0.7
- Microaveraged precision: 100/120 = .83

51

---

## Development Test Sets and Cross-validation

Training set     Development/tuning/held-out Set     Test Set

Metric: P/R/F1  or Accuracy

Cross-validation over multiple splits
- Handle sampling errors from different datasets
- Pool results over each split
- Compute pooled dev set performance

Training  Dev Test

Training Set     Dev Test

Dev Test     Training Set

Test Set

52