

Introduction

Regular  
Expression

Notations  
Examples

Text  
Classification

Naive Bayes  
Example

Language  
Modelling

Unigram, Bigram  
and N-gram

# Natural Language Processing

Nov 19, 2019

Introduction

Regular  
Expression

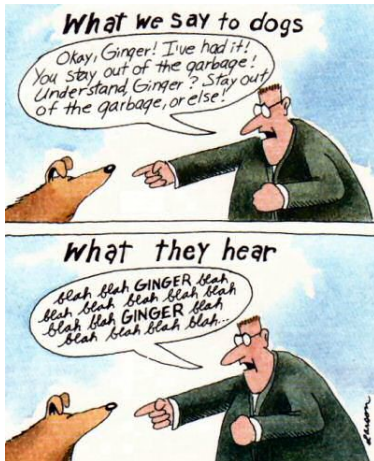
Notations  
Examples

Text  
Classification

Naive Bayes  
Example

Language  
Modelling

Unigram, Bigram  
and N-gram



# Who wrote the Federalist Papers?

Formal  
Language  
Processing

## Introduction

Regular  
Expression  
Notations  
Examples

Text  
Classification  
Naive Bayes  
Example

Language  
Modelling  
Unigram, Bigram  
and N-gram

- 1787-8: anonymous essays try to convince New York to ratify U.S Constitution: Jay, Madison, Hamilton.
- Authorship of 12 of the letters in dispute.
- 1963: solved by Mosteller and Wallace using Bayesian methods.

# Who wrote the Federalist Papers?

Formal  
Language  
Processing

## Introduction

Regular  
Expression  
Notations  
Examples

Text  
Classification  
Naive Bayes  
Example

Language  
Modelling  
Unigram, Bigram  
and N-gram

- 1787-8: anonymous essays try to convince New York to ratify U.S Constitution: Jay, Madison, Hamilton.
- Authorship of 12 of the letters in dispute.
- 1963: solved by Mosteller and Wallace using Bayesian methods.

By the end of this lecture we will see how to do that.

# What makes it hard?

## Introduction

### Regular Expression

Notations  
Examples

### Text Classification

Naive Bayes  
Example

### Language Modelling

Unigram, Bigram  
and N-gram

- Formal languages are:
  - *unambiguous*
- Natural languages are
  - *ambiguous*:
    - “He saw her duck” .
    - “Time flies like an arrow. Fruit flies like a banana”

# By the end of the class

## Introduction

### Regular Expression

Notations  
Examples

### Text Classification

Naive Bayes  
Example

### Language Modelling

Unigram, Bigram  
and N-gram

By the end of the class we will see how to do:

- 1 Text Classification. E.g. Spam detection, Authorship identification.
- 2 Spell Correction. E.g. Auto-correct.
- 3 Word suggestion.

# Regular Expressions

Natural  
Language  
Processing

Introduction

**Regular  
Expression**

Notations  
Examples

Text  
Classification

Naive Bayes  
Example

Language  
Modelling

Unigram, Bigram  
and N-gram

A formal language for specifying text strings.

# Notations

- Disjunctions [] :

Pattern	Matches
$[Ww]oodchuck$	woodchuck, Woodchuck
$[0123456789]$	Any single digit

- Disjunctions |:

Pattern	Matches
$abc def$	Find 'abc' or 'def'.
$a b ab$	Find 'a' or 'b' or 'ab'. Example: 'abc'



# Notations

- Ranges:

Pattern	Matches
$[A - Z]$	An uppercase letter.
$[a - z]$	A lowercase letter.
$[0 - 9]$	A single digit.

- Negation  $\wedge$ . (Note: Carat means negation only when its first in  $[]$ )

Pattern	Matches
$[\wedge A - Z]$	Not upper case
$[\wedge Ss]$	Not 'S' nor 's'
$[\wedge e\wedge]$	Not 'e' nor '^'
$a\wedge b$	Search for the pattern 'a^b'

# Notations (? \* . + ^ \$)

?	0 or 1 of previous character
*	0 or more of previous character
+	1 or more of previous character
.	Any character
^	Start anchor
\$	End anchor
\	Escape character

# Examples

- Pattern:  $\text{^[A-Z]}$ .

Which of them are matches? “Class”, “cSCE”, “420”.

# Examples

- Pattern:  $\^[A - Z]$ .

Which of them are matches? "Class", "cSCE", "420".

Class.

# Examples

- Pattern:  $\^[A - Z]$ .

Which of them are matches? "Class", "cSCE", "420".

Class.

- Pattern:  $\^[^A - Z]$ .

Which of them are matches? "Class", "cSCE", "420".

# Examples

- Pattern:  $\^[A - Z]$ .

Which of them are matches? "Class", "cSCE", "420".

Class.

- Pattern:  $\^[^A - Z]$ .

Which of them are matches? "Class", "cSCE", "420".

cSCE, 420.

# Examples

- Pattern:  $^[A - Z]$ .

Which of them are matches? "Class", "cSCE", "420".

Class.

- Pattern:  $^[^A - Z]$ .

Which of them are matches? "Class", "cSCE", "420".

cSCE, 420.

- Pattern:  $.\$$

Which of them are matches? "end", "end?", "end!", "end.".

# Examples

- Pattern:  $^[A-Z]$ .

Which of them are matches? "Class", "cSCE", "420".

Class.

- Pattern:  $^[^A-Z]$ .

Which of them are matches? "Class", "cSCE", "420".

cSCE, 420.

- Pattern:  $.\$$

Which of them are matches? "end", "end?", "end!", "end.".

end. end? end! end.



# Examples

- Pattern:  $^[A-Z]$ .

Which of them are matches? "Class", "cSCE", "420".

Class.

- Pattern:  $^[^A-Z]$ .

Which of them are matches? "Class", "cSCE", "420".

cSCE, 420.

- Pattern:  $.\$$

Which of them are matches? "end", "end?", "end!", "end.".

end. end? end! end.

- Pattern:  $\backslash.\$$

Which of them are matches? "end", "end?", "end!", "end.".

# Examples

- Pattern:  $^[A-Z]$ .

Which of them are matches? "Class", "cSCE", "420".

Class.

- Pattern:  $^[^A-Z]$ .

Which of them are matches? "Class", "cSCE", "420".

cSCE, 420.

- Pattern:  $.\$$

Which of them are matches? "end", "end?", "end!", "end.".

end. end? end! end.

- Pattern:  $\backslash.\$$

Which of them are matches? "end", "end?", "end!", "end.".

end..

# Examples

- Pattern: *colou?r*.

Which of them are matches? “color”, “colour”, “colouur”.

# Examples

- Pattern: *colou?r*.

Which of them are matches? “color”, “colour”, “colouur”.  
**color, colour.**

# Examples

- Pattern: *colou?r*.

Which of them are matches? “color”, “colour”, “colouur”.  
*color, colour.*

- Pattern: *colou + r*.

Which of them are matches? “color”, “colour”, “colouur”.

# Examples

- Pattern: *colou?r*.

Which of them are matches? “color”, “colour”, “colouur”.  
*color, colour.*

- Pattern: *colou + r*.

Which of them are matches? “color”, “colour”, “colouur”.  
*colour, colouur.*

# Examples

- Pattern:  $colou?r$ .

Which of them are matches? “color”, “colour”, “colouur”.  
**color, colour.**

- Pattern:  $colou + r$ .

Which of them are matches? “color”, “colour”, “colouur”.  
**colour, colour.**

- Pattern:  $colou * r$ .

Which of them are matches? “color”, “colour”, “colouur”.

# Examples

- Pattern:  $colou?r$ .

Which of them are matches? “color”, “colour”, “colouur”.  
**color, colour.**

- Pattern:  $colou + r$ .

Which of them are matches? “color”, “colour”, “colouur”.  
**colour, colour.**

- Pattern:  $colou * r$ .

Which of them are matches? “color”, “colour”, “colouur”.  
**color, colour, colour.**



# Examples

- Pattern:  $colou?r$ .

Which of them are matches? “color”, “colour”, “colouur”.  
**color, colour.**

- Pattern:  $colou + r$ .

Which of them are matches? “color”, “colour”, “colouur”.  
**colour, colour.**

- Pattern:  $colou * r$ .

Which of them are matches? “color”, “colour”, “colouur”.  
**color, colour, colour.**

- Pattern:  $colou.r$ .

Which of them are matches? “color”, “colour”, “colouur”.

# Examples

- Pattern:  $colou?r$ .

Which of them are matches? “color”, “colour”, “colouur”.  
**color, colour.**

- Pattern:  $colou + r$ .

Which of them are matches? “color”, “colour”, “colouur”.  
**colour, colour.**

- Pattern:  $colou * r$ .

Which of them are matches? “color”, “colour”, “colouur”.  
**color, colour, colour.**

- Pattern:  $colou.r$ .

Which of them are matches? “color”, “colour”, “colouur”.  
**colour.**

# Example

We need to find instances of “the” in a text.

# Example

We need to find instances of “the” in a text.

*the*

# Example

We need to find instances of “the” in a text.

*the* × 'The'

# Example

We need to find instances of “the” in a text.

*the* × 'The'

[Tt]he

# Example

We need to find instances of “the” in a text.

*the* × ‘The’

[*Tt*]he × ‘Theology’

# Example

We need to find instances of “the” in a text.

*the* × ‘The’

[*Tt*]he × ‘Theology’

[ $\hat{A} - Z a - z$ ][*Tt*]he[ $\hat{A} - Z a - z$ ]



# Text classification

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

# Text Classification

- Inputs:
  - Document  $d$ .
  - Fixed set of classes  $C = \{c_1, c_2, \dots, c_n\}$ .
- Output:
  - A predicted class  $c \in C$

# Naive Bayes

Natural  
Language  
Processing

Introduction

Regular  
Expression

Notations  
Examples

Text  
Classification

**Naive Bayes**  
Example

Language  
Modelling

Unigram, Bigram  
and N-gram

Relies on simple representation of document – Bag of Words.

# Naive Bayes

Relies on simple representation of document – Bag of Words.

For a document  $d$  and a class  $c$

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

# Naive Bayes Classifier

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

MAP - Maximum a posteriori (most likely class).

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

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

# Naive Bayes Classifier

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

Let's say that the document is represented by  $n$  features  
 $x_1, x_2, \dots, x_n$

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n|c)P(c)$$

# Assumptions

Mathematical  
Language  
Processing

Introduction

Regular  
Expression

Notations  
Examples

Text  
Classification

Naive Bayes  
Example

Language  
Modelling

Unigram, Bigram  
and N-gram

**Bag of words:** Position of words does not matter.

**Conditional Independence:** The feature probabilities  $P(x_i|c)$  are independent given the class  $c$ .

$$P(x_1, x_2, \dots, x_n|c) = \prod_{i=1}^n P(x_i|c)$$

# Bag of word 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.



# Bag of word representation

great	2
love	2
recommend	1
laugh	1
happy	1
...	...

# Naive Bayes: Learning

Natural  
Language  
Processing

Introduction

Regular  
Expression

Notations  
Examples

Text  
Classification

Naive Bayes  
Example

Language  
Modelling

Unigram, Bigram  
and N-gram

What do we need?

Training set of  $m$  hand-labeled documents

$(d_1, c_1), \dots, (d_m, c_m)$

# Naive Bayes: Learning

Let  $N_D$  be the number of documents, and  $N_{c_j}$  be the number of documents present in class  $c_j$ .

Let  $V_{c_j}$  be the set of all words in the documents of class  $c_j$

Now we find the maximum likelihood estimates:

$$\hat{P}(c_j) = \frac{N_{c_j}}{N_D}$$

$$\hat{P}(w_i|c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V_{c_j}} \text{count}(w, c_j)}$$

Now we can classify a document  $d$  by:

$$c_d = \operatorname{argmax}_{c_j \in \mathcal{C}} \hat{P}(c_j) \prod_{w_i \in d} \hat{P}(w_i|c_j)$$

# Naive Bayes: Learning

Natural  
Language  
Processing

Introduction

Regular  
Expression

Notations  
Examples

Text  
Classification

**Naive Bayes**  
Example

Language  
Modelling

Unigram, Bigram  
and N-gram

What if we come across an unknown word in the document  $d$ .  
Let  $w_u$  be the unknown word  $\hat{P}(w_u|c_j) = 0, \forall c_j$ .

# Laplace smoothing

Let  $V$  be the set of all words in the test documents, i.e.,

$$V = \cup_{c_j} V_{c_j}$$

Add one word for the unknown word in the vocabulary.

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j) + 1}{\sum_{w \in V_{c_j}} \text{count}(w, c_j) + |V| + 1}$$

So, for all unknown words, we have:

$$\hat{P}(w_u | c_j) = \frac{1}{\sum_{w \in V_{c_j}} \text{count}(w, c_j) + |V| + 1}$$

# Example

Training set:

#	Text	Class
1	Carla Betty Carla	a
2	Carla Carla Suzanne	a
3	Carla Matt	a
4	Taylor Jessica Carla	b

# Example

Training set:

#	Text	Class
1	Carla Betty Carla	a
2	Carla Carla Suzanne	a
3	Carla Matt	a
4	Taylor Jessica Carla	b

$$\hat{P}(a) =$$

# Example

Training set:

#	Text	Class
1	Carla Betty Carla	a
2	Carla Carla Suzanne	a
3	Carla Matt	a
4	Taylor Jessica Carla	b

$$\hat{P}(a) = 3/4$$

$$\hat{P}(b) =$$



# Example

Training set:

#	Text	Class
1	Carla Betty Carla	a
2	Carla Carla Suzanne	a
3	Carla Matt	a
4	Taylor Jessica Carla	b

$$\hat{P}(a) = 3/4$$

$$\hat{P}(b) = 1/4$$

# Example

Training set:

#	Text	Class
1	Carla Betty Carla	a
2	Carla Carla Suzanne	a
3	Carla Matt	a
4	Taylor Jessica Carla	b

# Example

Training set:

#	Text	Class
1	Carla Betty Carla	a
2	Carla Carla Suzanne	a
3	Carla Matt	a
4	Taylor Jessica Carla	b

$$\hat{P}(\text{Carla}|a) =$$

# Example

Training set:

#	Text	Class
1	Carla Betty Carla	a
2	Carla Carla Suzanne	a
3	Carla Matt	a
4	Taylor Jessica Carla	b

$$\hat{P}(Carla|a) = (5 + 1)/(8 + 6 + 1) = 6/15$$
$$\hat{P}(Taylor|a) =$$

# Example

Training set:

#	Text	Class
1	Carla Betty Carla	a
2	Carla Carla Suzanne	a
3	Carla Matt	a
4	Taylor Jessica Carla	b

$$\hat{P}(\text{Carla}|a) = (5 + 1)/(8 + 6 + 1) = 6/15$$

$$\hat{P}(\text{Taylor}|a) = (0 + 1)/(8 + 6 + 1) = 1/15$$

# Example

Document  $d_5$ : *Carla Carla Carla Taylor Jessica*

# Example

Document  $d_5$ : *Carla Carla Carla Taylor Jessica*

$$\hat{P}(a) = 3/4 \text{ and } \hat{P}(b) = 1/4$$

$$\hat{P}(Carla|a) = 6/15, \hat{P}(Carla|b) = 2/10$$

$$\hat{P}(Taylor|a) = 1/15, \hat{P}(Taylor|b) = 2/10$$

$$\hat{P}(Jessica|a) = 1/15, \hat{P}(Jessica|b) = 2/10$$

# Example

Document  $d_5$ : *Carla Carla Carla Taylor Jessica*

$$\hat{P}(a) = 3/4 \text{ and } \hat{P}(b) = 1/4$$

$$\hat{P}(Carla|a) = 6/15, \hat{P}(Carla|b) = 2/10$$

$$\hat{P}(Taylor|a) = 1/15, \hat{P}(Taylor|b) = 2/10$$

$$\hat{P}(Jessica|a) = 1/15, \hat{P}(Jessica|b) = 2/10$$

$$P(a|d_5) =$$



# Example

Document  $d_5$ : *Carla Carla Carla Taylor Jessica*

$$\hat{P}(a) = 3/4 \text{ and } \hat{P}(b) = 1/4$$

$$\hat{P}(Carla|a) = 6/15, \hat{P}(Carla|b) = 2/10$$

$$\hat{P}(Taylor|a) = 1/15, \hat{P}(Taylor|b) = 2/10$$

$$\hat{P}(Jessica|a) = 1/15, \hat{P}(Jessica|b) = 2/10$$

$$P(a|d_5) = 3/4 \times (6/15)^3 \times 1/15 \times 1/15 \approx 0.0002$$

$$P(b|d_5) =$$

# Example

Document  $d_5$ : *Carla Carla Carla Taylor Jessica*

$$\hat{P}(a) = 3/4 \text{ and } \hat{P}(b) = 1/4$$

$$\hat{P}(Carla|a) = 6/15, \hat{P}(Carla|b) = 2/10$$

$$\hat{P}(Taylor|a) = 1/15, \hat{P}(Taylor|b) = 2/10$$

$$\hat{P}(Jessica|a) = 1/15, \hat{P}(Jessica|b) = 2/10$$

$$P(a|d_5) = 3/4 \times (6/15)^3 \times 1/15 \times 1/15 \approx 0.0002$$

$$P(b|d_5) = 1/4 \times (2/10)^3 \times 2/10 \times 2/10 \approx 0.00008$$

# Naive Bayes

Natural  
Language  
Processing

Introduction

Regular  
Expression

Notations  
Examples

Text  
Classification

Naive Bayes  
Example

Language  
Modelling

Unigram, Bigram  
and N-gram

Naive Bayes is not so naive!!

- Robust to Irrelevant Features.
- Optimal if the independence assumptions hold.
- A good dependable baseline for text classification. - **There exists other classifiers that give better accuracy**

# Federalist Papers

Natural  
Language  
Processing

Introduction

Regular  
Expression

Notations  
Examples

Text  
Classification

Naive Bayes

**Example**

Language  
Modelling

Unigram, Bigram  
and N-gram

Discussion: Federalist papers.  
E.g. What training set do we need?

# Language Modelling

Natural  
Language  
Processing

Introduction

Regular  
Expression

Notations  
Examples

Text  
Classification

Naive Bayes  
Example

**Language  
Modelling**

Unigram, Bigram  
and N-gram

**Goal:** Assign probability to a sentence.

# Language Modelling

Natural  
Language  
Processing

Introduction

Regular  
Expression

Notations  
Examples

Text  
Classification

Naive Bayes  
Example

**Language  
Modelling**

Unigram, Bigram  
and N-gram

**Goal:** Assign probability to a sentence.  
**Why?**

# Language Modelling

Natural  
Language  
Processing

Introduction

Regular  
Expression

Notations  
Examples

Text  
Classification

Naive Bayes  
Example

**Language  
Modelling**

Unigram, Bigram  
and N-gram

**Goal:** Assign probability to a sentence.

**Why?**

- **Machine Translation.**  $P(\text{high winds tonight}) > P(\text{large winds tonight})$

# Language Modelling

Formal  
Language  
Processing

Introduction

Regular  
Expression

Notations  
Examples

Text  
Classification

Naive Bayes  
Example

Language  
Modelling

Unigram, Bigram  
and N-gram

**Goal:** Assign probability to a sentence.

**Why?**

- **Machine Translation.**  $P(\text{high winds tonight}) > P(\text{large winds tonight})$
- **Spell Correction.**  $P(\text{about fifteen minutes from}) > P(\text{about fifteen minuets from})$



# Language Modelling

Formal  
Language  
Processing

Introduction

Regular  
Expression

Notations  
Examples

Text  
Classification

Naive Bayes  
Example

Language  
Modelling

Unigram, Bigram  
and N-gram

**Goal:** Assign probability to a sentence.

**Why?**

- **Machine Translation.**  $P(\text{high winds tonight}) > P(\text{large winds tonight})$
- **Spell Correction.**  $P(\text{about fifteen minutes from}) > P(\text{about fifteen minuets from})$
- **Speech Recognition.**  $P(\text{I saw a van}) > P(\text{eyes awe of an})$

# Language Modelling

- **Goal:** compute the probability of a sentence or sequence of words.  $P(W) = P(w_1, w_2, \dots, w_n)$
- **Related task:** probability of an upcoming word.  
 $P(w_i | w_1, w_2, \dots, w_{i-1})$

# Language Modelling

- **Goal:** compute the probability of a sentence or sequence of words.  $P(W) = P(w_1, w_2, \dots, w_n)$
- **Related task:** probability of an upcoming word.  
 $P(w_i | w_1, w_2, \dots, w_{i-1})$

## Chain rule:

$$\begin{aligned} & P(x_1, x_2, \dots, x_n) \\ &= P(x_1)P(x_2|x_1)P(x_3|x_1, x_2) \cdots P(x_n|x_1, x_2, \dots, x_{n-1}) \\ &= \prod_i P(x_i|x_1, x_2, \dots, x_{i-1}) \end{aligned}$$

# Example

$$P(\text{its water is so transparent that}) = P(\text{its}) \times P(\text{water}|\text{its}) \times P(\text{so}|\text{its, water}) \times P(\text{transparent}|\text{its, water, is, so}) \times P(\text{that}|\text{its, water, is, so, transparent})$$

Can we count?

$$P(\text{that}|\text{its, water, is, so, transparent}) = \frac{P(\text{its, water, is, so, transparent, that})}{P(\text{its, water, is, so, transparent})}$$

# Example

$$P(\text{its water is so transparent that}) = P(\text{its}) \times P(\text{water}|\text{its}) \times P(\text{so}|\text{its, water}) \times P(\text{transparent}|\text{its, water, is, so}) \times P(\text{that}|\text{its, water, is, so, transparent})$$

Can we count?

$$\begin{aligned} & P(\text{that}|\text{its, water, is, so, transparent}) \\ = & \frac{P(\text{its, water, is, so, transparent, that})}{P(\text{its, water, is, so, transparent})} \end{aligned}$$

No. Too many possibilities.

# Markov Assumption

Take only the  $k$  words preceding it.

$$P(w_i | w_1, w_2, \dots, w_{i-1}) \approx P(w_i | w_{i-k} \dots, w_{i-1})$$

# Markov Assumption

Take only the  $k$  words preceding it.

$$P(w_i | w_1, w_2, \dots, w_{i-1}) \approx P(w_i | w_{i-k} \dots, w_{i-1})$$

$$P(\text{that} | \text{its, water, is, so, transparent}) = P(\text{that} | \text{transparent})$$

or,

$$P(\text{that} | \text{its, water, is, so, transparent}) = P(\text{that} | \text{so, transparent})$$

# Unigram, Bigram and N-gram

## Unigram model:

$$P(w_1, w_2, \dots, w_{n-1}, w_n) \approx \prod_i P(w_i)$$

## Bigram model:

$$P(w_i | w_1, w_2, \dots, w_{i-1}) \approx P(w_i | w_{i-1})$$

## N-gram model:

Extension to trigram, 4-gram, 5-gram, etc.



# Discussion

Natural  
Language  
Processing

Introduction

Regular  
Expression

Notations  
Examples

Text  
Classification

Naive Bayes  
Example

Language  
Modelling

Unigram, Bigram  
and N-gram

## 1. Spell correction

# Discussion

Natural  
Language  
Processing

Introduction

Regular  
Expression

Notations  
Examples

Text  
Classification

Naive Bayes  
Example

Language  
Modelling

Unigram, Bigram  
and N-gram

1. Spell correction
2. Word suggestion.