

# COMP90042 Workshop Week 3

## Text Classification & N-gram Language Model

Zenan Zhai

The University of Melbourne

16 March 2014

# Table of Contents

Text Classifications

N-gram Language Model

# Table of Contents

Text Classifications

N-gram Language Model

# Definition

Think of an application of text classification.

# Definition

Think of an application of text classification.

- Sentiment analysis
- Author identification
- Fact checking
- ...

What makes text classification challenging?

# Definition

Think of an application of text classification.

- Sentiment analysis
- Author identification
- Fact checking
- ...

What makes text classification challenging?

- How to learn document representation?
- How to do feature selection?
- How to deal with data sparsity?

# Text classification Models

- K-Nearest Neighbors (KNN)
- Decision Tree
- Naive Bayes
- Logistic Regression
- Support Vector Machine
- ...

# K-Nearest Neighbors

Vote by the label of  $K$  nearest instances in the training set

Similarity Measure

- Euclidean distance  $d(A, B) = \sqrt{\sum (a_i - b_i)^2}$



# K-Nearest Neighbors

Vote by the label of  $K$  nearest instances in the training set

Similarity Measure

- Euclidean distance  $d(A, B) = \sqrt{\sum (a_i - b_i)^2}$

Not ideal, the measure is greatly affected by **document length**.

# K-Nearest Neighbors

Vote by the label of  $K$  nearest instances in the training set

Similarity Measure

- Cosine similarity  $\cos(A, B) = \frac{A \cdot B}{\|A\| \|B\|}$

# K-Nearest Neighbors

Vote by the label of  $K$  nearest instances in the training set

Similarity Measure

- Cosine similarity  $\cos(A, B) = \frac{A \cdot B}{\|A\| \|B\|}$

Better than Eculidean distance, but suffers from  
**curse of dimensionality**

# Decision Tree

Entropy

$$H(Y) = - \sum_{y \in Y} p(y) \log(p(y))$$

Conditional Entropy

$$H(Y|X) = \sum_{x \in X} p(x) H(Y|X = x)$$

Information Gain

$$IG(Y|a) = H(Y) - H(Y|a)$$

# Decision Tree

Entropy

$$H(Y) = - \sum_{y \in Y} p(y) \log(p(y))$$

Conditional Entropy

$$H(Y|X) = \sum_{x \in X} p(x) H(Y|X = x)$$

Information Gain

$$IG(Y|a) = H(Y) - H(Y|a)$$

Tends to prefer **rare features** which might only appears in a few documents.

# Naive Bayes

Every feature is assumed to be independent

$$\begin{aligned} P(y|x_1, x_2, x_3, \dots, x_n) &\propto P(y, x_1, x_2, x_3, \dots, x_n) \\ &= P(y) \prod_{i=1}^n P(x_i|y) \end{aligned}$$

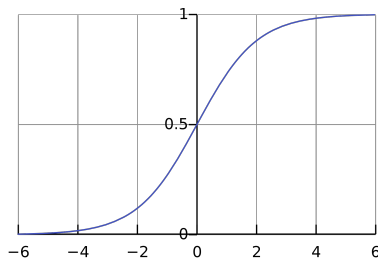
# Naive Bayes

Every feature is assumed to be independent

$$\begin{aligned} P(y|x_1, x_2, x_3, \dots, x_n) &\propto P(y, x_1, x_2, x_3, \dots, x_n) \\ &= P(y) \prod_{i=1}^n P(x_i|y) \end{aligned}$$

Still work on dataset **large set of features**, but suffers from biases caused by **uninformative features**.

# Logistic Regression

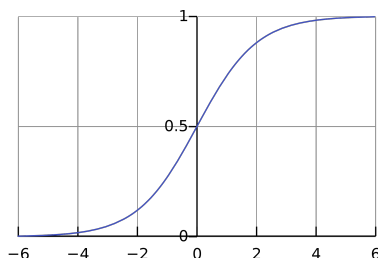


Put linear combination of features in logistic function

$$P(y) = \sigma(y) = \frac{1}{1 + \exp(-WX + b)}$$



# Logistic Regression

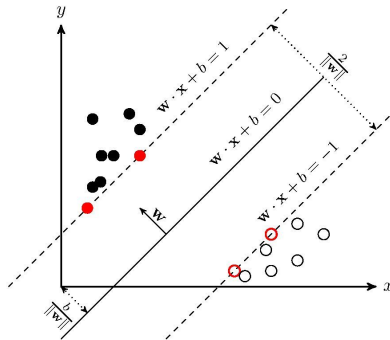


Put linear combination of features in logistic function

$$P(y) = \sigma(y) = \frac{1}{1 + \exp(-WX + b)}$$

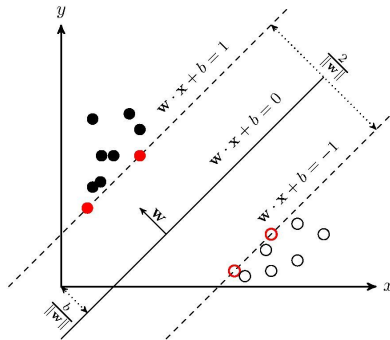
**Relaxed the independent assumption** and solves the problem caused by uninformative features by taking weighted sum.

# Support Vector Machine



Select the decision boundary which maximize the distance to the support vectors.

# Support Vector Machine



Select the decision boundary which maximize the distance to the support vectors.

An very effective methods, but no natural support  
**multi-classification**

# Table of Contents

Text Classifications

N-gram Language Model

# Uni-gram model

## Corpus

1. how much wood would a wood chuck chuck if a wood chuck would chuck wood
2. a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood

# Uni-gram model

## Corpus

1. how much wood would a wood chuck chuck if a wood chuck would chuck wood
2. a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood

## Word counts

a	chuck	could	he	how	if	much	the	wood	would	</s>	Total
4	9	1	1	1	2	1	1	8	4	2	34

# Uni-gram model

## Corpus

1. how much wood would a wood chuck chuck if a wood chuck would chuck wood
2. a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood

## Word counts

a	chuck	could	he	how	if	much	the	wood	would	</s>	Total
4	9	1	1	1	2	1	1	8	4	2	34

Why is <s> left out?

# Uni-gram model

## Word counts

a	chuck	could	he	how	if	much	the	wood	would	</s>	Total
4	9	1	1	1	2	1	1	8	4	2	34

## Sentences

A: a wood could chuck

B: wood would a chuck



# Uni-gram model

## Word counts

a	chuck	could	he	how	if	much	the	wood	would	</s>	Total
4	9	1	1	1	2	1	1	8	4	2	34

## Sentences

A: a wood could chuck

B: wood would a chuck

$$\begin{aligned}P(A) &= P(a)P(\text{wood})P(\text{could})P(\text{chuck})P(</s>) \\&= \frac{4}{34} \times \frac{8}{34} \times \frac{1}{34} \times \frac{9}{34} \times \frac{2}{34} \approx 1.27 \times 10^{-5} \\P(B) &= P(\text{wood})P(\text{would})P(a)P(\text{chuck})P(</s>) \\&= \frac{8}{34} \times \frac{4}{34} \times \frac{4}{34} \times \frac{9}{34} \times \frac{2}{34} \approx 5.07 \times 10^{-5}\end{aligned}$$

# Uni-gram model with Laplacian smoothing

## Word counts

a	chuck	could	he	how	if	much	the	wood	would	</s>	Total
4	9	1	1	1	2	1	1	8	4	2	34

## Sentences

A: a wood could chuck

B: wood would a chuck

$$\begin{aligned}P_L(A) &= P_L(a)P_L(\text{wood})P_L(\text{could})P_L(\text{chuck})P_L(</s>) \\&= \frac{5}{45} \times \frac{9}{45} \times \frac{2}{45} \times \frac{10}{45} \times \frac{3}{45} \approx 1.46 \times 10^{-5} \\P_L(B) &= P_L(\text{wood})P_L(\text{would})P_L(a)P_L(\text{chuck})P_L(</s>) \\&= \frac{9}{45} \times \frac{5}{45} \times \frac{5}{45} \times \frac{10}{45} \times \frac{3}{45} \approx 3.66 \times 10^{-5}\end{aligned}$$

# Bi-gram model

## Corpus

1. how much wood would a wood chuck chuck if a wood chuck would chuck wood
2. a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood

## Sentences

A: a wood could chuck

B: wood would a chuck

# Bi-gram model

## Corpus

1. how much wood would a wood chuck chuck if a wood chuck would chuck wood
2. a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood

## Sentences

A: a wood could chuck

B: wood would a chuck

$$\begin{aligned}P(A) &= P(a|<s>)P(\text{wood}|a)P(\text{could}|\text{wood})P(\text{chuck}|\text{could})P(</s>|\text{chuck}) \\&= \frac{1}{2} \times \frac{4}{4} \times \frac{0}{8} \times \frac{1}{1} \times \frac{0}{9} = 0\end{aligned}$$

$$\begin{aligned}P(B) &= P(\text{wood}|<s>)P(\text{would}|\text{wood})P(a|\text{would})P(\text{chuck}|a)P(</s>|\text{chuck}) \\&= \frac{0}{2} \times \frac{1}{8} \times \frac{1}{4} \times \frac{0}{4} \times \frac{0}{9} = 0\end{aligned}$$

# Bi-gram model with Laplacian smoothing

## Corpus

1. how much wood would a wood chuck chuck if a wood chuck would chuck wood
2. a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood

## Sentences

A: a wood could chuck

B: wood would a chuck

$$\begin{aligned}P_L(A) &= P_L(a|<s>)P_L(\text{wood}|a)P_L(\text{could}|\text{wood})P_L(\text{chuck}|\text{could})P_L(</s>|\text{chuck}) \\&= \frac{2}{13} \times \frac{5}{15} \times \frac{1}{19} \times \frac{2}{12} \times \frac{1}{20} \approx 2.25 \times 10^{-5} \\P_L(B) &= P_L(\text{wood}|<s>)P_L(\text{would}|\text{wood})P_L(a|\text{would})P_L(\text{chuck}|a)P_L(</s>|\text{chuck}) \\&= \frac{1}{13} \times \frac{2}{19} \times \frac{2}{15} \times \frac{1}{15} \times \frac{1}{20} \approx 3.60 \times 10^{-6}\end{aligned}$$

# Tri-gram model

## Corpus

1. how much wood would a wood chuck chuck if a wood chuck would chuck wood
2. a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood

## Sentences

A: a wood could chuck

B: wood would a chuck

# Tri-gram model

## Corpus

1. how much wood would a wood chuck chuck if a wood chuck would chuck wood
2. a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood

## Sentences

A: a wood could chuck

B: wood would a chuck

$$\begin{aligned}P(A) &= P(a|\langle s \rangle \langle s \rangle)P(\text{wood}|\langle s \rangle a)\dots P(\langle /s \rangle | \text{could chuck}) \\ &= \frac{1}{2} \times \frac{1}{1} \times \frac{0}{4} \times \frac{0}{0} \times \frac{0}{1} = ?\end{aligned}$$

$$\begin{aligned}P(B) &= P(\text{wood}|\langle s \rangle \langle s \rangle)P(\text{would}|\langle s \rangle \text{wood})\dots P(\langle /s \rangle | a \text{ chuck}) \\ &= \frac{0}{2} \times \frac{0}{0} \times \frac{1}{1} \times \frac{0}{1} \times \frac{0}{0} = ?\end{aligned}$$

# Tri-gram model with Laplacian smoothing

## Corpus

1. how much wood would a wood chuck chuck if a wood chuck would chuck wood
2. a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood

## Sentences

A: a wood could chuck

B: wood would a chuck

$$\begin{aligned}P_L(A) &= P_L(a|<s> <s>)P_L(\text{wood}|<s> a) \cdots P_L(</s>|could chuck) \\&= \frac{2}{13} \times \frac{2}{12} \times \frac{1}{15} \times \frac{1}{11} \times \frac{1}{12} \approx 1.30 \times 10^{-5}\end{aligned}$$

$$\begin{aligned}P_L(B) &= P_L(\text{wood}|<s> <s>)P_L(\text{would}|<s> wood) \cdots P_L(</s>|a chuck) \\&= \frac{1}{13} \times \frac{1}{11} \times \frac{2}{12} \times \frac{1}{12} \times \frac{1}{11} \approx 8.83 \times 10^{-6}\end{aligned}$$



# Kneser-Ney Smoothing

What is continuation count?

# Kneser-Ney Smoothing

What is continuation count?

- number of word types in the vocabulary which appears before a word  $w$
- $|\{w_{i-1} : C(w_{i-1}, w_i) > 0\}|$

What is Kneser-Ney smoothing?

# Kneser-Ney Smoothing

What is continuation count?

- number of word types in the vocabulary which appears before a word  $w$
- $|\{w_{i-1} : C(w_{i-1}, w_i) > 0\}|$

What is Kneser-Ney smoothing?

- use continuation probability instead of trivial uni-gram prob
- can be used in either back-off and interpolation

$$P_{cont}(w_i) = \frac{|\{w_{i-1} : C(w_{i-1}, w_i) > 0\}|}{\sum_{w_i \in V} |\{w_{i-1} : C(w_{i-1}, w_i) > 0\}|}$$

# KN Smoothing - Example

## Corpus

1. how much wood would a wood chuck chuck if a wood chuck would chuck wood
2. a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood

## Continuation counts

a =  
he =  
if =  
the =  
</s> =

could =  
how =  
much =  
would =

# KN Smoothing - Example

## Corpus

1. how much wood would a wood chuck chuck if a wood chuck would chuck wood
2. a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood

## Continuation counts

a = 2  
he = 1  
if = 1  
the = 1  
</s> = 1

could = 1  
how = 0  
much = 1  
would = 2

# KN Smoothing - Example

## Corpus

1. how much wood would a wood chuck chuck if a wood chuck would chuck wood
2. a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood

## Continuation counts

chuck=

a = 2

he = 1

if = 1

the = 1

</s> = 1

wood=

could = 1

how = 0

much = 1

would = 2

# KN Smoothing - Example

## Corpus

1. how much wood would a wood chuck chuck if a wood chuck would chuck wood
2. a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood

## Continuation counts

chuck= 4

a = 2

he = 1

if = 1

the = 1

</s> = 1

wood= 4

could = 1

how = 0

much = 1

would = 2

# KN Smoothing - Example

## Continuation counts

chuck = 4

a = 2

he = 1

if = 1

the = 1

</s> = 1

wood = 4

could = 1

how = 0

much = 1

would = 2

## Continuation probabilities

$$P_{cont}(\text{chuck})$$

$$P_{cont}(\text{wood})$$



# KN Smoothing - Example

## Continuation counts

chuck = 4

a = 2

he = 1

if = 1

the = 1

</s> = 1

wood = 4

could = 1

how = 0

much = 1

would = 2

## Continuation probabilities

$$\begin{aligned} P_{cont}(\text{chuck}) &= \frac{\#_{cont}(\text{chuck})}{\#_{cont}(\text{a}) + \dots + \#_{cont}(\text{</s>}) + \#_{cont}(\text{chuck}) + \#_{cont}(\text{wood})} \\ &= \frac{4}{2 + 1 + 1 + 0 + 1 + 1 + 1 + 2 + 1 + 4 + 4} \\ P_{cont}(\text{wood}) \end{aligned}$$

# KN Smoothing - Example

## Continuation counts

chuck = 4

a = 2

he = 1

if = 1

the = 1

</s> = 1

wood = 4

could = 1

how = 0

much = 1

would = 2

## Continuation probabilities

$$\begin{aligned} P_{cont}(\text{chuck}) &= \frac{\#_{cont}(\text{chuck})}{\#_{cont}(\text{a}) + \dots + \#_{cont}(\text{</s>}) + \#_{cont}(\text{chuck}) + \#_{cont}(\text{wood})} \\ &= \frac{4}{2 + 1 + 1 + 0 + 1 + 1 + 1 + 2 + 1 + 4 + 4} \end{aligned}$$

$$\begin{aligned} P_{cont}(\text{wood}) &= \frac{\#_{cont}(\text{wood})}{\#_{cont}(\text{a}) + \dots + \#_{cont}(\text{</s>}) + \#_{cont}(\text{chuck}) + \#_{cont}(\text{wood})} \\ &= \frac{4}{2 + 1 + 1 + 0 + 1 + 1 + 1 + 2 + 1 + 4 + 4} \end{aligned}$$

# Back-off and Interpolation

## Back-off

- Use lower-order  $n$ -gram model if higher-order is unseen.

$$P(w_i|w_{i-1}) = \begin{cases} \frac{c(w_i, w_{i-1}) - D}{c(w_{i-1})}, & \text{if } c(w_i, w_{i-1}) > 0 \\ \alpha(w_{i-1}) \times \frac{P(w_i)}{\sum_{w_j: C(w_{i-1}, w_j)=0} P(w_j)}, & \text{otherwise} \end{cases}$$

## Interpolation

- Take weighted average sum of all orders

$$P(w_i|w_{i-1}) = \lambda P(w_i|w_{i-1}) + (1 - \lambda)P(w_i)$$

# Evaluation

Recall the objective of language model:

- Modeling probability for an arbitrary sequence of  $m$  words.

Evaluate based on probability of all sequences in test set

$$PP(w_1, w_2, w_3, \dots, w_m) = \sqrt[m]{\frac{1}{P(w_1, w_2, w_3, \dots, w_m)}}$$

- Inverted prob. : lower perplexity  $\rightarrow$  better model
- Normalization : take  $m^{th}$  root of sequence prob. ,  
 $m = length(S)$

# Take aways

- Text classification
  - ▶ Applications
  - ▶ Models (Pros & Cons)
- $n$ -gram language model (calculation)
- Smoothing
  - ▶ problems they address
  - ▶ Pros & Cons