

# Homework One: Natural Language Processing

Quin'darius Lyles-Woods

September 8, 2021

## 1 Regular Expressions

### 1.1 Binary Strings

$(0|1)^*$

### 1.2 Email

$^{\wedge}\backslash S+@\backslash S+.\backslash S+^{\$}$

### 1.3 Integers

$\backslash d$

### 1.4 Phone Number

$^{\wedge}(\backslash +\backslash d\{1,2\}\backslash s)?\backslash (? \backslash d\{3\}\backslash )?[\backslash s.-]\backslash d\{3}[\backslash s.-]\backslash d\{4}^{\$}$

## 2 Tokens and Vocabulary

Text:

“The quick brown fox jumps over the lazy dog.”

Tokens 9

Types 8

### 3 Text Normalization

#### Steps of Normalization:

**Tokenization** Turning text into words.

**Example:** My name is quin and I live in Georgia.

This contains nine tokens or words.

**Lemmatization** Reducing a collection of words to their root meaning

**Example:** My name is quin and I am living in Georgia. I lived in Georgia for a while.

Living and lived will be reduced to live.

**Stemming** Close to Lemmatization, with Stemming you're taking away the suffixes.

**Example:** My name is quin and I am living in Georgia. I lived in Georgia for a while.

Living and lived will be reduced to liv.

**Sentence Segmentation** Breaking up the text with the delimiters of . , ! or ?

**Example:** My name is quin and I am living in Georgia. I lived in Georgia for a while. I like living here.

**Sentence One** My name is quin and I am living in Georgia.

**Sentence Two** I lived in Georgia for a while.

**Sentence Three** I like living here.

**Spelling Correction** Checking for spelling mistakes.

**Example:** My name is quin afd I live in Georgia.

The word afd should be and.

**Non-Standard Words** Phone numbers, emails, dates, etc.

**Example:** My name is quin and my phone number is 704-470-7036, and email is qlyleswo@student.kennesaw.edu.

**Phone Number** 704-470-7036

**Email** qlyleswo@student.kennesaw.edu

### 4 Similarity Distance with Edit Distance Algorithm

**String One:** Spokesman confirms

**String Two:** Spokeswoman said

## 4.1 Steps for Distance Matrix

You must compare the strings and see what needs to be replaced to get them to equal each other.

**The operations are:**

Insert +1

Delete +1

Replace +2

The total after comparing the strings will be 13. The steps are assumed to be shown sufficiently in the distance matrices.

## 4.2 Distance Matrix

		S	p	o	k	e	w	o	m	a	n		s	a	i	d
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
S	1	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
p	2	1	0	1	2	3	4	5	6	7	8	9	10	11	12	13
o	3	2	1	0	1	2	3	4	5	6	7	8	9	10	11	12
k	4	3	2	1	0	1	2	3	4	5	6	7	8	9	10	11
e	5	4	3	2	1	0	1	2	3	4	5	6	7	8	9	10
s	6	5	4	3	2	1	2	3	4	5	6	7	6	7	8	9
m	7	6	5	4	3	2	3	4	3	4	5	6	7	8	9	10
a	8	7	6	5	4	3	4	5	4	3	4	5	6	7	8	9
n	9	8	7	6	5	4	5	6	5	4	3	4	5	6	7	8
	10	9	8	7	6	5	6	7	6	5	4	3	4	5	6	7
c	11	10	9	8	7	6	7	8	7	6	5	4	5	6	7	8
o	12	11	10	9	8	7	8	7	8	7	6	5	6	7	8	9
n	13	12	11	10	9	8	9	8	9	8	7	6	7	8	9	10
f	14	13	12	11	10	9	10	9	10	9	8	7	8	9	10	11
i	15	14	13	12	11	10	11	10	11	10	9	8	9	10	9	10
r	16	15	14	13	12	11	12	11	12	11	10	9	10	11	10	11
m	17	16	15	14	13	12	13	12	11	12	11	10	11	12	11	12
s	18	17	16	15	14	13	14	13	12	13	12	11	10	11	12	13

### 4.3 Backtracing Matrix Check

		S	p	o	k	e	w	o	m	a	n	s	a	i	d	
S p o k e s m a n  c o n f i r m s	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
	1	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
	2	1	0	1	2	3	4	5	6	7	8	9	10	11	12	13
	3	2	1	0	1	2	3	4	5	6	7	8	9	10	11	12
	4	3	2	1	0	1	2	3	4	5	6	7	8	9	10	11
	5	4	3	2	1	0	1	2	3	4	5	6	7	8	9	10
	6	5	4	3	2	1	2	3	4	5	6	7	6	7	8	9
	7	6	5	4	3	2	3	4	3	4	5	6	7	8	9	10
	8	7	6	5	4	3	4	5	4	3	4	5	6	7	8	9
	9	8	7	6	5	4	5	6	5	4	3	4	5	6	7	8
	10	9	8	7	6	5	6	7	6	5	4	3	4	5	6	7
	11	10	9	8	7	6	7	8	7	6	5	4	5	6	7	8
	12	11	10	9	8	7	8	7	8	7	6	5	6	7	8	9
	13	12	11	10	9	8	9	8	9	8	7	6	7	8	9	10
	14	13	12	11	10	9	10	9	10	9	8	7	8	9	10	11
	15	14	13	12	11	10	11	10	11	10	9	8	9	10	9	10
	16	15	14	13	12	11	12	11	12	11	10	9	10	11	10	11
	17	16	15	14	13	12	13	12	11	12	11	10	11	12	11	12
18	17	16	15	14	13	14	13	12	13	12	11	10	11	12	13	

## 5 Language Model

Text:

The day was grey and bitter cold, and the dogs would not take the scent. The big black bitch had taken one sniff at the bear tracks, backed off, and skulked back to the pack with her tail between her legs.

### 5.1 Unigram Model

the 6/30

and 3/30

her 2/30

. 2/30

, 2/30

## 5.2 Bigram Model

The	day	was	grey	and
day	0	1	0	0
was	1	0	1	0
grey	0	1	0	1
and	0	0	1	0

Choose to do a smaller table and show the iterations because it wouldn't fit all in one page. but its just 0's and 1's for the whole thing there is nothing with 2 to my knowledge for the **raw bigram count**.

- the day
- day was
- was grey
- grey and
- and bitter
- bitter cold,
- cold, and
- and the
- the dogs
- dogs would
- would not
- not take
- take the
- the scent.
- scent. the
- the big
- big black
- black bitch
- bitch had
- had taken
- taken one
- one sniff
- sniff at
- at the
- the bear
- bear tracks,
- tracks, backed
- backed off,
- off, and
- and skulked
- skulked back
- back to
- to the
- the pack
- pack with
- with her
- her tail
- tail between
- between her
- her legs.

## 6 Unigram Perplexity

You are given a training set of 30 numbers that consists of 21 zeros and 1 each of the other digits 1-9. Now we see the following test set: 0 0 0 0 0 3 0 0 0 0. What is the unigram perplexity?

0	1	2	3	4	5	6	7	8	9
21	1	1	1	1	1	1	1	1	1

Table 1: Normalizing by Unigrams

0	1	2	3	4	5	6	7	8	9
.7	.03	.03	.03	.03	.03	.03	.03	.03	.03

Table 2: Computing Probabilities

$$\text{Unigram Perplexity} = .7 * (.03 * .03 * .03 * .03 * .03 * .03 * .03 * .03 * .03) = 2.1 \quad (1)$$