# English Word Suggestion Based on Part of Speech Ngram

Hamana Laboratory, Gunma University

Borann Chanrathnak

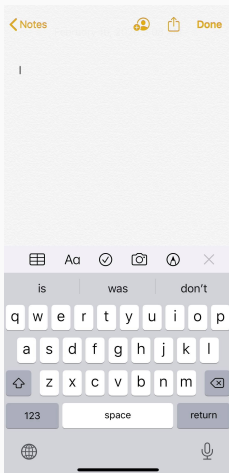February 17, 2020

## Table of Contents

# Motivation

## Motivation

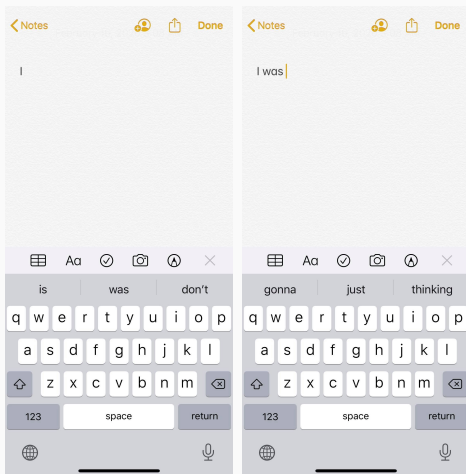- Word suggestion on mobile devices

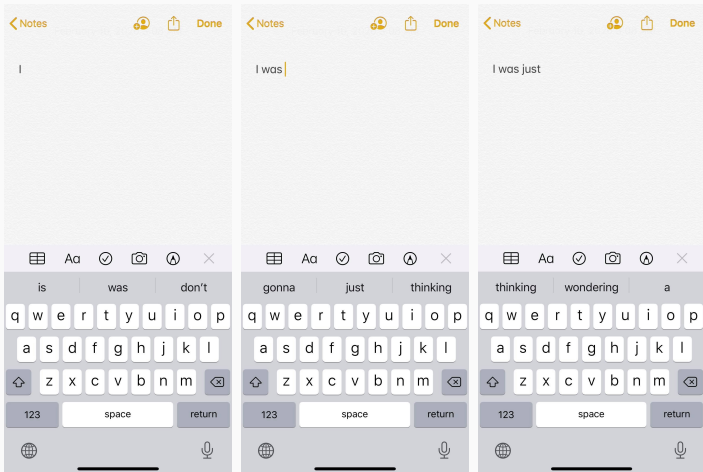# Motivation

- Word suggestion on mobile devices

# Motivation

- Word suggestion on mobile devices

# Motivation

- Word suggestion on mobile devices

# Part Of Speech Ngram (POS-Ngram)

## What does POS-Ngram mean?

**N-gram** is a contiguous sequence of n items from a given sample of text or speech

**Ex:**  I study english

- **unigram (1-gram)** : (I,), (study,), (english,)
- **bigram (2-gram)** : (I, study), (study, english)
- **trigram (3-gram)** : (I, study, english)

## What does POS-Ngram mean?

**N-gram** is a contiguous sequence of n items from a given sample of text or speech

**Ex:** I study english

- **unigram (1-gram)** : (I,), (study,), (english,)
- **bigram (2-gram)** : (I, study), (study, english)
- **trigram (3-gram)** : (I, study, english)

**N-gram model** is one of statistical langauge models for predicting the next item (word) based on Markov assumption, and usually abbreviated as N-gram.

## What does POS-Ngram mean?

**N-gram** is a contiguous sequence of n items from a given sample of text or speech

**Ex:** I study english

- **unigram (1-gram)** : (I,), (study,), (english,)
- **bigram (2-gram)** : (I, study), (study, english)
- **trigram (3-gram)** : (I, study, english)

**N-gram model** is one of statistical langauge models for predicting the next item (word) based on Markov assumption, and usually abbreviated as N-gram.

**POS-Ngram** is an improved model using part of speech as a class indicator.

**How To Compute Probability of a Sentence?**

How can we compute the joint probability of a sentence?
*Ex:* "an apple is on the table"

$$P(an\ apple\ is\ on\ the\ table) = ?$$

## Notation

## Notation

- To represent the probability of a particular random variable $X_i$ taking on the value "the", or $P(X_i = "the")$ we will use simplification $P(the)$.

### Notation

- To represent the probability of a particular random variable $X_i$ taking on the value "the", or $P(X_i = "the")$ we will use simplification $P(the)$.

- We represent a sequence of N words either as $w_1 \cdots w_n$ or $w_1^n$

### Notation

- To represent the probability of a particular random variable $X_i$ taking on the value "the", or $P(X_i = "the")$ we will use simplification $P(the)$.

- We represent a sequence of N words either as $w_1 \cdots w_n$ or $w_1^n$

- **Ex:** "an apple is on the table"
  $\rightarrow w_1 = "an", w_2 = "apple", w_3 = "is", \cdots$

## Notation

- To represent the probability of a particular random variable $X_i$ taking on the value "the", or $P(X_i = "the")$ we will use simplification $P(the)$.

- We represent a sequence of N words either as $w_1 \cdots w_n$ or $w_1^n$

- **Ex:** "an apple is on the table"
  $\rightarrow w_1 = "an", \ w_2 = "apple", \ w_3 = "is", \cdots$
  $\rightarrow w_1^4 = "an \ apple \ is \ on"$

## Chain Rule Of Probability

**In General**

$$P(X_1 \cdots X_n) = P(X_1)P(X_2 \mid X_1)P(X_3 \mid X_1^2) \cdots P(X_n \mid X_1^{n-1})$$

# Chain Rule Of Probability

## In General

$$P(X_1 \cdots X_n) = P(X_1)P(X_2 \mid X_1)P(X_3 \mid X_1^2) \cdots P(X_n \mid X_1^{n-1})$$

## By applying chain rule to a sequence of words

$$P(w_1 \cdots w_n) = P(w_1)P(w_2 \mid w_1)P(w_3 \mid w_1^2) \cdots P(w_n \mid w_1^{n-1})$$

## Chain Rule Of Probability

### In General

$$P(X_1 \cdots X_n) = P(X_1)P(X_2 \mid X_1)P(X_3 \mid X_1^2) \cdots P(X_n \mid X_1^{n-1})$$

### By applying chain rule to a sequence of words

$$P(w_1 \cdots w_n) = P(w_1)P(w_2 \mid w_1)P(w_3 \mid w_1^2) \cdots P(w_n \mid w_1^{n-1})$$

### Examples

- $P(an\ apple) = P(an) \times P(apple \mid an)$
- $P(an\ apple\ is) = P(an) \times P(apple \mid an) \times P(is \mid an\ apple)$

## Let's Apply Chain Rule

By applying chain rule to the phrase "*an apple is on the*":

$$P("\text{an apple is on the}") = P(an) \times P(apple \mid an) \times P(is \mid an\ apple)$$
$$\times P(on \mid an\ apple\ is) \times P(the \mid an\ apple\ is\ on)$$

**In practice chain rule does not help**

- We don't know the way to compute the exact probability of a word given a long sequence of preceding words, $P(w_n \mid w_1^{n-1})$

## In practice chain rule does not help

- We don't know the way to compute the exact probability of a word given a long sequence of preceding words, $P(w_n \mid w_1^{n-1})$
- Language is creative, and any particular context might have never occured before

## The Presence of N-gram

The general equation for this n-gram approximation to the conditional probability of the next word in a sequence is

$$P(w_n \mid w_1^{n-1}) \approx P(w_n \mid w_{n-N+1}^{n-1})$$

## The Presence of N-gram

The general equation for this n-gram approximation to the conditional probability of the next word in a sequence is

$$P(w_n \mid w_1^{n-1}) \approx P(w_n \mid w_{n-N+1}^{n-1})$$

**In case of Bigram (2-gram)**
When we use bigram model to predict the conditional probability of the next word, we can approximate by

$$P(w_n \mid w_1^{n-1}) \approx P(w_n \mid w_{n-1})$$

## The Presence of N-gram

The general equation for this n-gram approximation to the conditional probability of the next word in a sequence is

$$P(w_n \mid w_1^{n-1}) \approx P(w_n \mid w_{n-N+1}^{n-1})$$

**In case of Bigram (2-gram)**
When we use bigram model to predict the conditional probability of the next word, we can approximate by

$$P(w_n \mid w_1^{n-1}) \approx P(w_n \mid w_{n-1})$$

**Example**

$$P(the \mid an\ apple\ is\ on) \approx P(the \mid on)$$

## Bigram in Practice

By supposing that C is the counts of word from a corpus.

$$P(w_n \mid w_{n-1}) = \frac{C(w_{n-1}w_n)}{\sum_{w \in dict} C(w_{n-1}w)}$$

## Bigram in Practice

By supposing that C is the counts of word from a corpus.

$$P(w_n \mid w_{n-1}) = \frac{C(w_{n-1}w_n)}{\sum_{w \in dict} C(w_{n-1}w)}$$

**Example**

$$P(the \mid on) = \frac{C(on\ the)}{\sum_x C(on\ x)}$$

## The probability of a complete word sequence

Given the **Bigram** assumption for the probability of an individual word, we can compute the probability of complete word sequence by:

$$P(w_1^n) \approx \prod_{k=1}^{n} P(w_k \mid w_{k-1})$$

**Chain Rule**

$$P("an\ apple\ is\ on\ the") = P(an) \times P(apple \mid an) \times P(is \mid an\ apple)$$
$$\times P(on \mid an\ apple\ is) \times P(the \mid an\ apple\ is\ on)$$

## The probability of a complete word sequence

Given the **Bigram** assumption for the probability of an individual word, we can compute the probability of complete word sequence by:

$$P(w_1^n) \approx \prod_{k=1}^{n} P(w_k \mid w_{k-1})$$

**Bigram Assumption**

$$P("\text{an apple is on the}") = P(\text{apple} \mid \text{an}) \times P(\text{is} \mid \text{apple})$$
$$\times P(\text{on} \mid \text{is}) \times P(\text{the} \mid \text{on})$$

## OOV & Smoothing

**OOV** stands for **Out Of Vocabulary**

- Words appear only in a *test set* but not in the *training set*.
- OOV problem occurs even when we work on big data.

- Laplace smoothing (Add-one smoothing)

- Add-k smoothing

- Interpolation

- $\cdots$

**Interpolation**

We mix the probability estimates from all the n-gram estimators, weighing and combining the trigram, bigram, and unigram counts.

## Interpolation

**In case of Bigram**

$$\hat{P}(table \mid the) = \lambda_1 P(table \mid the) + \lambda_2 P(table)$$

## Interpolation

**In case of Bigram**

$$\hat{P}(table \mid the) = \lambda_1 P(table \mid the) + \lambda_2 P(table)$$

**In case of Trigram**

$$\hat{P}(table \mid on\ the) = \lambda_1 P(table \mid on\ the)$$
$$+ \lambda_2 P(table \mid the)$$
$$+ \lambda_3 P(table)$$

## Interpolation

**In case of Bigram**

$$\hat{P}(table \mid the) = \lambda_1 P(table \mid the) + \lambda_2 P(table)$$

**In case of Trigram**

$$\hat{P}(table \mid on\ the) = \lambda_1 P(table \mid on\ the)$$
$$+ \lambda_2 P(table \mid the)$$
$$+ \lambda_3 P(table)$$

where we choose $\lambda_i$ such that $\sum_i \lambda_i = 1$

**How are the $\lambda$ values set?**

- They can be learned from a **held-out** corpus
- Can be found by **EM** algorithm
- For the purpose of this project, we assume without loss of generality that $\lambda_i > \lambda_j (\forall i < j)$

## POS-Ngram

**Formula**

$$P(w_n \mid w_1^{n-1}) = \sum_{c_n} P(w_n \mid c_n) \times P(c_n \mid c_{n-N+1}^{n-1})$$

where $c_n$ is a class of $w_n$

(In this context, $c_n$ is considered a part of speech of the word $w_n$)

## POS-Ngram

**Formula**

$$P(w_n \mid w_1^{n-1}) = \sum_{c_n} P(w_n \mid c_n) \times P(c_n \mid c_{n-N+1}^{n-1})$$

where $c_n$ is a class of $w_n$

(In this context, $c_n$ is considered a part of speech of the word $w_n$)

**Example**

| $c_1$:Noun | cat, dog, thought, $\cdots$ |
|---|---|
| $c_2$:Verb | go, speak, $\cdots$ |
| $c_1$:Noun, $c_2$:Verb | play, $\cdots$ |

## POS-Ngram

### Formula

$$P(w_n \mid w_1^{n-1}) = \sum_{c_n} P(w_n \mid c_n) \times P(c_n \mid c_{n-N+1}^{n-1})$$

where $c_n$ is a class of $w_n$

(In this context, $c_n$ is considered a part of speech of the word $w_n$)

### Example

| $c_1$:Noun | cat, dog, thought, $\cdots$ |
|---|---|
| $c_2$:Verb | go, speak, $\cdots$ |
| $c_1$:Noun, $c_2$:Verb | play, $\cdots$ |

It means that to predict the next word $w_n$

1. Compute $P(c_n \mid c_{n-N+1}^{n-1})$
2. $P(w_n \mid c_n) = \frac{C(w_n, c_n)}{C(c_n)}$
3. $w_n^*$ is defined by $\arg \max_{w_n} P(w_n \mid w_1^{n-1})$

## POS-Ngram

**Formula**

$$P(w_n \mid w_1^{n-1}) = \sum_{c_n} P(w_n \mid c_n) \times P(c_n \mid c_{n-N+1}^{n-1})$$

where $c_n$ is a class of $w_n$

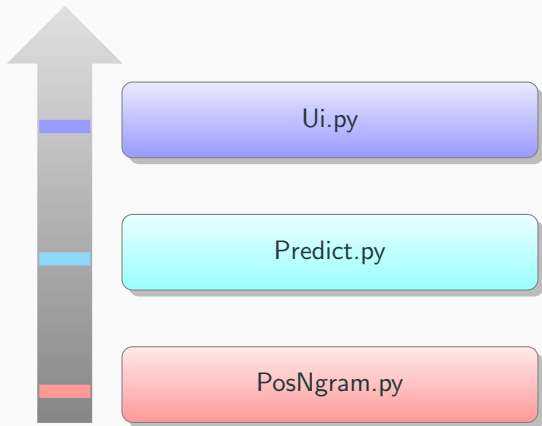(In this context, $c_n$ is considered a part of speech of the word $w_n$)

**Example**

| $c_1$:Noun | cat, dog, thought, $\cdots$ |
|---|---|
| $c_2$:Verb | go, speak, $\cdots$ |
| $c_1$:Noun, $c_2$:Verb | play, $\cdots$ |

It means that to predict the next word $w_n$

1. Compute $P(c_n \mid c_{n-N+1}^{n-1}) \sim P(w_n \mid w_{n-N+1}^{n-1})$
2. $P(w_n \mid c_n) = \frac{C(w_n, c_n)}{C(c_n)}$
3. $w_n^*$ is defined by $\arg\max_{w_n} P(w_n \mid w_1^{n-1})$

# Implementation

## Implementation

- Programming Language : Python
- Ui : Tkinter
- NLTK (Natural Language ToolKit)

## References

1. Speech and Language Processing (Chapter 4) (Daniel Jurafsky - Stanford University, 1999)
2. N-gram Language Modeling Tutorial (Lecture notes courtesy of Prof. Mari Ostendorf, 2007-06-21)
3. Probabilistic Language Model (Chapter 3) (Kenji KITA, University of Tokyo Press, 1999)

DEMO

Thank You.