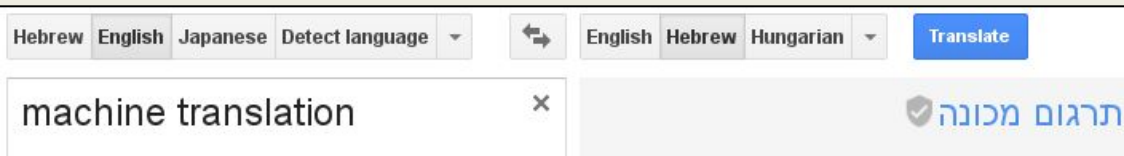# Language Models

Human Language from a Computational Perspective
April 11, 2018

# Natural Language Processing

## Algorithms that understand or generate human language



Automatic summarization is the process of reducing a text document with a computer program in order to create a summary that retains the most important points of the original document. Technologies that can make a coherent summary take into account variables such as length, writing style and syntax. Automatic data summarization is part of machine learning and data mining. The main idea of summarization is to find a representative subset of the data, which contains the *information* of the entire set. Summarization technologies are used in a large number of sectors in industry today. An example of the use of summarization technology is search engines such as Google. Other examples include document summarization, image collection summarization and video summarization. Document summarization, tries to automatically create a *representative summary* or *abstract* of the entire document, by finding the most *informative* sentences. Similarly, in image summarization the system finds the most representative and important (or salient) images. Similarly, in consumer videos one would want to remove the boring or repetitive scenes, and extract out a much shorter and concise version of the video.

Automatic summarization: reducing text with a computer to retain the most important points.

Hebrew | English | Japanese | Detect language | English | Hebrew | Hungarian | **Translate**

machine translation ✕    תרגום מכונה

what is question answering?

**Question Answering** (QA) is a computer science discipline within the fields of information retrieval and natural language processing (NLP), which is concerned with building systems that automatically **answer questions** posed by humans in a natural language.

# Statistical Language Model

How likely is each of these sentences?

PLEASE MAKE ME A CUP OF COFFEE
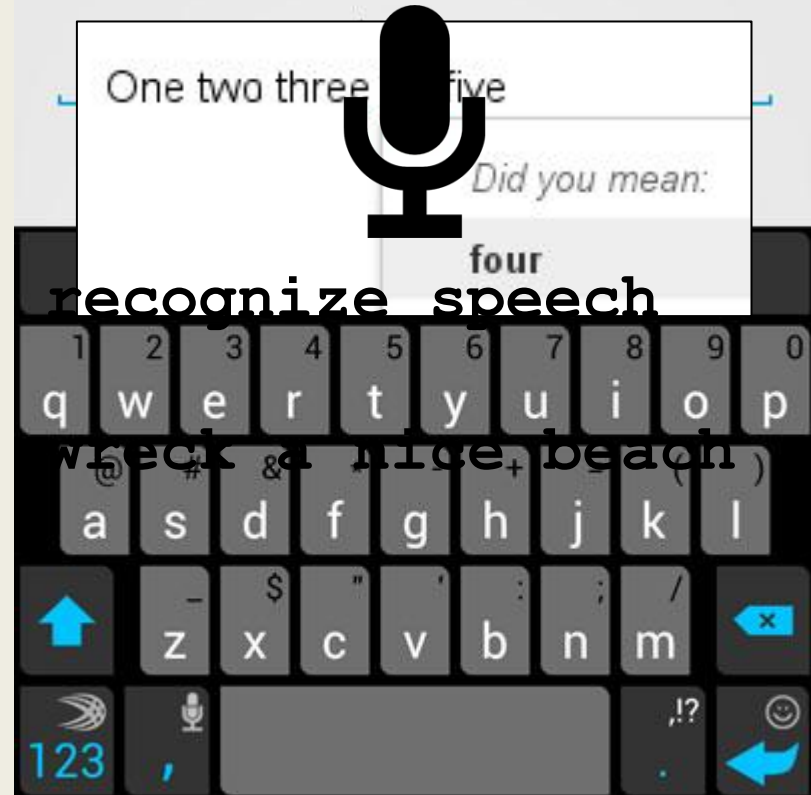
PLEASE MAKE ME A CUP OF BUTTER

PLEASE MAKE ME A CUP OF BOTTLE

PLEASE MAKE ME A CUP OF DREAM

PLEASE MAKE ME A CUP OF PLEASE

# Uses of Language Models

- Typing prediction

- Spelling correction

- Speech recognition

- Many more

# Algorithm

Instructions for manipulating data.

Can get parameters as **input**.

Returns an **output**.

input →  Algorithm  → output

# Pseudocode

Notation to describe algorithms.

Not a programming language, but clear enough for humans.

# Algorithm to find maximum

Find the largest number in a list.

$$[3, 1, 4, 16, 0, 2] \rightarrow 16$$

$$[1, 2, 1, 1, 1] \rightarrow 2$$

$$[-3, -2, 0, -1] \rightarrow 0$$

# Algorithm to find maximum

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

    i ← i + 1

  return m

▷ L is a list of numbers

▷ assign the first number to m

▷ assign 2 to i

▷ repeat while i is at most len(L)

▷ the i'th number is larger than m

▷ assign the i'th number to m

▷ increase i by 1

▷ output is the value of m

# Algorithm to find maximum

```
max(L):
    m ← L[1]
    i ← 2
    while i ≤ len(L):
        if L[i] > m:
            m ← L[i]
        i ← i + 1
    return m
```

**Comments**

▷ L is a list of numbers

▷ assign the first number to m

▷ assign 2 to i

▷ repeat while i is at most len(L)

▷ the i'th number is larger than m

▷ assign the i'th number to m

▷ increase i by 1

▷ output is the value of m

# Algorithm to find maximum

**max**(L):

▷ L is a list of numbers

m ← L[1]   ▷ assign the first number to m

i ← 2   ▷ assign 2 to i

while i ≤ len(L):   ▷ repeat while i is at most len(L)

    if L[i] > m:   ▷ the i'th number is larger than m

        m ← L[i]   ▷ assign the i'th number to m

    i ← i + 1   ▷ increase i by 1

return m   ▷ output is the value of m

# Algorithm to find maximum

**max(L)**    **Parameters**    ▷ L is a list of numbers

    m ← L[1]                 ▷ assign the first number to m

    i ← 2                    ▷ assign 2 to i

    while i ≤ len(L):        ▷ repeat while i is at most len(L)

       if L[i] > m:          ▷ the i'th number is larger than m

         m ← L[i]         ▷ assign the i'th number to m

A function can get more than one parameter, but **max** gets just one

# Algorithm to find maximum

**max**(L):

   m ← L[1]

   i ← 2  **Variables**

   while i ≤ len(L):

      if L[i] > m:

         m ← L[i]

      i ← i + 1

   return m

▷ L is a list of numbers

▷ assign the first number to m

▷ assign 2 to i

▷ repeat while i is at most len(L)

▷ the i'th number is larger than m

▷ assign the i'th number to m

▷ increase i by 1

▷ output is the value of m

# Algorithm to find maximum

**max**(L): ▷ L is a list of numbers

   m ← L[1] ▷ assign the first number to m

   i ← 2   **Assignment** ▷ assign 2 to i

   while i ≤ len(L): ▷ repeat while i is at most len(L)

      if L[i] > m: ▷ the i'th number is larger than m

        m ← L[i] ▷ assign the i'th number to m

      i ← i + 1 ▷ increase i by 1

   return m ▷ output is the value of m

# Algorithm to find maximum

**max**(L):                                   ▹ L is a list of numbers

    m ← L[1] **Indexing**       ▹ assign the first number to m

    i ← 2                        ▹ assign 2 to i

    while i ≤ len(L):            ▹ repeat while i is at most len(L)

        if L[i] > m:             ▹ the i'th number is larger than m

            m ← L[i]          ▹ assign the i'th number to m

        i ← i + 1                ▹ increase i by 1

    return m                     ▹ output is the value of m

# Algorithm to find maximum

**max**(L):                              ▷ L is a list of numbers

    m ← L[1] **Function**     ▷ assign the first number to m

    i ← 2    **call**  ▷ assign 2 to i

    while i ≤ len(L):          ▷ repeat while i is at most len(L)

      if L[i] > m:           ▷ the i'th number is larger than m

       m ← L[i]           ▷ assign the i'th number to m

The function **len** returns the number of elements (length) of a list

# Algorithm to find maximum

**max**(L):

   m ← L[1]

   i ← 2   **Loop**

while i ≤ len(L):

   if L[i] > m:

      m ← L[i]

   i ← i + 1

return m

▷ L is a list of numbers

▷ assign the first number to m

▷ assign 2 to i

▷ repeat while i is at most len(L)

▷ the i'th number is larger than m

▷ assign the i'th number to m

▷ increase i by 1

▷ output is the value of m

# Algorithm to find maximum

**max**(L):

    m ← L[1]

    i ← 2   **Condition**

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

▷ L is a list of numbers

▷ assign the first number to m

▷ assign 2 to i

▷ repeat while i is at most len(L)

▷ the i'th number is larger than m

▷ assign the i'th number to m

▷ increase i by 1

▷ output is the value of m

# Algorithm to find maximum

**max**(L):                              ▷ L is a list of numbers

    m ← L[1]                      ▷ assign the first number to m

    i ← 2                          ▷ assign 2 to i

    while i ≤ len(L):              ▷ repeat while i is at most len(L)

        if L[i] > m:               ▷ the i'th number is larger than m

            m ← L[i]           ▷ assign the i'th number to m

**Output** i ← i + 1                     ▷ increase i by 1

    return m                       ▷ output is the value of m

# Algorithm to find maximum

**max**(L): **Indentation**      ▷ L is a list of numbers

m ← L[1]      ▷ assign the first number to m

i ← 2      ▷ assign 2 to i

while i ≤ len(L):      ▷ repeat while i is at most len(L)

    if L[i] > m:      ▷ the i'th number is larger than m

        m ← L[i]      ▷ assign the i'th number to m

i ← i + 1      ▷ increase i by 1

return m      ▷ output is the value of m

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m

i

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m

i

# Running the algorithm

**max**(L):

   m ← L[1]

   i ← 2

   while i ≤ len(L):

      if L[i] > m:

         m ← L[i]

      i ← i + 1

   return m

L = [3, 1, 4, 16, 0, 2]

m = 3

i

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m = 3

i = 2

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m = 3

i = 2

    len(L) = 6

    2 ≤ 6 ✔

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m = 3

i = 2

L[i] = L[2] = 1

1 > 3 ✘

# Running the algorithm

**max**(L):
    m ← L[1]
    i ← 2
    while i ≤ len(L):
        if L[i] > m:
            m ← L[i]
        i ← i + 1
    return m

L = [3, 1, 4, 16, 0, 2]
m = 3
i = 3

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m = 3

i = 3

len(L) = 6

3 ≤ 6 ✔

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m = 3

i = 3

L[i] = L[3] = 4

4 > 3 ✔

# Running the algorithm

**max**(L):
    m ← L[1]
    i ← 2
    while i ≤ len(L):
        if L[i] > m:
            m ← L[i]
        i ← i + 1
    return m

L = [3, 1, 4, 16, 0, 2]
m = 4
i = 3

L[i] = L[3] = 4

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

    i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m = 4

i = 4

# Running the algorithm

**max**(L):

   m ← L[1]

   i ← 2

while i ≤ len(L):

    if L[i] > m:

      m ← L[i]

    i ← i + 1

return m

L = [3, 1, 4, 16, 0, 2]

m = 4

i = 4

len(L) = 6

4 ≤ 6 ✔

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m = 4

i = 4

    L[i] = L[4] = 16

    16 > 4 ✔

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m = 16

i = 4

L[i] = L[4] = 16

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m = 16

i = 5

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m = 16

i = 5

len(L) = 6

5 ≤ 6 ✔

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m = 16

i = 5

L[i] = L[5] = 0

0 > 16 ✗

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m = 16

i = 6

# Running the algorithm

**max**(L):
    m ← L[1]
    i ← 2
    while i ≤ len(L):
        if L[i] > m:
            m ← L[i]
        i ← i + 1
    return m

L = [3, 1, 4, 16, 0, 2]
m = 16
i = 6

len(L) = 6
6 ≤ 6 ✔

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m = 16

i = 6

L[i] = L[6] = 2

2 > 16 ✗

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

    return m

L = [3, 1, 4, 16, 0, 2]

m = 16

i = 7

# Running the algorithm

**max**(L):

   m ← L[1]

   i ← 2

  while i ≤ len(L):

     if L[i] > m:

       m ← L[i]

     i ← i + 1

  return m

L = [3, 1, 4, 16, 0, 2]

m = 16

i = 7

   len(L) = 6

   7 ≤ 6  ✗

# Running the algorithm

**max**(L):

    m ← L[1]

    i ← 2

    while i ≤ len(L):

        if L[i] > m:

            m ← L[i]

        i ← i + 1

return m

L = [3, 1, 4, 16, 0, 2]

m = 16

i = 7

output: 16

# Finding index of maximum

**Index** of largest number in a list.

$$[3, 1, 4, 16, 0, 2] \rightarrow 4$$

$$[1, 2, 1, 1, 1] \rightarrow 2$$

$$[-3, -2, 0, -1] \rightarrow 3$$

# Finding index of maximum

**argmax**(L):

    a ← 1

    i ← 2

    while i ≤ len(L):

        if L[i] > L[a]:

            a ← i

        i ← i + 1

    return a

▷ L is a list of numbers

▷ index of the first element

▷ index of the second element

▷ repeat while i is at most len(L)

▷ i'th number is larger than a'th

▷ assign i to a

▷ increase i by 1

▷ output: index of largest number

# Back to language models

Given a list of tokens, predict the most

likely token to follow.

PLEASE MAKE ME A CUP OF <u>TEA</u>

# Tokenization

We represent a string:

“I’M LATE!”, HE SAID.

As a **list** of tokens:

| “ | I | ‘M | LATE | ! | “ | , | HE | SAID | . |
|---|---|----|------|---|---|---|----|------|---|

# Language model algorithm

Predict the next token in the list.

PLEASE MAKE ME A CUP OF → COFFEE

ONE TWO THREE → FOUR

WHAT IS YOUR PHONE → NUMBER

But the list of tokens is not enough.

We also need to know the language.

# Corpora

A **text corpus** is used to analyze the distribution of words.

# Counts table

To represent token counts, we map strings to numbers.

Some ready-made counts:

books.google.com/ngrams

| | |
|---|---|
| , | 775 |
| THE | 630 |
| . | 392 |
| " | 345 |
| AND | 339 |
| A | 337 |
| TO | 277 |

# Algorithm to count words

Count all words in a tokenized corpus and return a table of counts.

[I, AM, SAM, .,

SAM, I, AM.,

I, DO, NOT, LIKE,

GREEN, EGGS, AND, HAM, .]

$\rightarrow$

| I | 3 | LIKE | 1 |
|---|---|------|---|
| AM | 2 | GREEN | 1 |
| SAM | 2 | EGGS | 1 |
| . | 3 | AND | 1 |
| DO | 1 | HAM | 1 |
| NOT | 1 | | |

# Algorithm to count words

**count**(L):              ▹ L is a list of tokens

    C1 ← [0]          ▹ create a table of zeros

    i ← 1               ▹ assign 1 to i

    while i ≤ len(L):    ▹ repeat while i is at most len(L)

        t ← L[i]        ▹ get token at position i

        C1[t] ← C1[t] + 1  ▹ increase count for t by 1

        i ← i + 1      ▹ increase i by 1

    return C1        ▹ output is the counts table

# Algorithm to count words

**count**(L):

C1 ← [0]

i ← 1

while i ≤ len(L):

t ← L[i]

C1[t] ← C1[t] + 1

i ← i + 1

return C1

**Using a word as an index to a table**

▷ L is a list of tokens

▷ create a table of zeros

▷ assign 1 to i

▷ repeat while i is at most len(L)

▷ get token at position i

▷ increase count for t by 1

▷ increase i by 1

▷ output is the counts table

# Word counts

Example counts from *Alice's Adventures in Wonderland* (1866) by Lewis Carroll

| | |
|---|---|
| , | 775 |
| THE | 630 |
| . | 392 |
| " | 345 |
| AND | 339 |
| A | 337 |
| TO | 277 |

# Unigram Language Model

Easiest: always predict

the most frequent token:

I WISH I **,**

C1 =

(**Unigram** counts)

C1[,] = 775

C1[THE] = 630

| | |
|---|---|
| , | **775** |
| THE | 630 |
| . | 392 |
| " | 345 |
| AND | 339 |
| A | 337 |
| TO | 277 |

# Bigram counts

We can also count

**bigrams** (pairs of words)

$$C2 =$$

| | | |
|---|---|---|
| , | THE | 530 |
| AND | THE | 320 |
| . | " | 89 |
| SHE | SAID | 65 |
| | ... | |
| I | 'M | 20 |
| I | DO | 10 |

$C2[\text{AND, THE}] = 320$

$C2[\text{SHE, SAID}] = 65$

# Bigram Language Model

Look only at the **last** token to predict the next:

I WISH **I 'M**

C2[I, ·] =

(**Bigram** counts starting with I)

| I | 'M | **20** |
|---|---|---|
| I | DO | 10 |
| I | 'LL | 10 |
| I | 'VE | 10 |
| I | SHOULD | 8 |
| I | MUST | 7 |
| I | THINK | 7 |

C2[I, 'M] = 20

C2[I, DO] = 10

# Trigram Language Model

Look at the **two** last tokens to predict the next:

I WISH I COULD

$$C3[\text{WISH, I, } \cdot] = \begin{array}{|ll|} \hline \text{WISH I COULD} & \textbf{20} \\ \text{WISH I HAD} & 10 \\ \hline \end{array}$$

(**Trigram** counts starting with WISH I)

C3[WISH, I, COULD] = 20

C3[WISH, I, HAD] = 10

# *n*-gram Language Model

Look at the ***n − 1*** last tokens to predict the next:

<mark>I WISH I COULD</mark>

C4[I, WISH, I, ·] =

(**4-gram** counts starting with I WISH I):

C3[I, WISH, I, COULD] = 2

C3[I, WISH, I, HAD] = 2

| | |
|---|---|
| **I WISH I COULD** | **2** |
| I WISH I HAD | 2 |

# Algorithm to count *n*-grams

**count**(L, n):                          ▷ L: list of tokens, n: a number

    C ← [0]                          ▷ create a table of zeros

    i ← 1                            ▷ assign 1 to i

    while i ≤ len(L) − n + 1:        ▷ repeat while i is at most len(L) − n + 1

        T ← L[i, ..., i + n − 1]     ▷ get n tokens starting at i

        C[T] ← C[T] + 1              ▷ increase count for T by 1

        i ← i + 1                    ▷ increase i by 1

    return C                         ▷ output is the counts table

# Algorithm to count *n*-grams

**count**(L, n):

C ← [0]

i ← 1

while i ≤ len(L) − n + 1:

T ← L[i, ..., i + n − 1]

C[T] ← C[T] + 1

i ← i + 1

return C

## Getting several elements from a list

▹ L: list of tokens, n: a number

▹ create a table of zeros

▹ assign 1 to i

▹ repeat while i is at most len(L) − n + 1

▹ get n tokens starting at i

▹ increase count for T by 1

▹ increase i by 1

▹ output is the counts table

# Algorithm to count *n*-grams

**count**(L, n):

C ← [0]

i ← 1

while i ≤ len(L) − n + 1:

    T ← L[i, …, i + n − 1]

    C[T] ← C[T] + 1

    i ← i + 1

return C

**Using an *n*-gram as an index to a table**

▷ L: list of tokens, n: a number

▷ create a table of zeros

▷ assign 1 to i

▷ repeat while i is at most len(L) − n + 1

▷ get n tokens starting at i

▷ increase count for T by 1

▷ increase i by 1

▷ output is the counts table

# Unigram algorithm

**unigram**(L, C1):

    return argmax(C1)

▷ L: tokens, C1: unigram counts

▷ token with highest count

# Unigram algorithm

**unigram**(L, C1):

    return argmax(C1)

**Function call**

▷ L: tokens, C1: unigram counts

▷ token with highest count

# Unigram algorithm

**unigram**(L, C1):
    return argmax(C1)

▷ L: tokens, C1: unigram counts

▷ token with highest count

Ignores L and always predicts the same word...

# Bigram algorithm

**bigram**(L, C2):

    k ← len(L)

    t ← L[k]

    return argmax(C2[t, ·])

▷ L: tokens, C2: bigram counts

▷ length of L

▷ last token in L

    ▷ bigram with highest count,

▷ among the bigrams starting with t

# Bigram algorithm

**bigram**(L, C2):  ▷ L: tokens, C2: bigram counts

    k ← len(L)  ▷ length of L

    t ← L[k]  ▷ last token in L

    return argmax(C2[t, ·])  ▷ bigram with highest count,

▷ among the bigrams starting with t

**Getting part of the table**

# Trigram algorithm

**trigram**(L, C3):
    k ← len(L)
    T ← L[k − 1, k]
    return argmax(C3[T, ·])

▷ L: tokens, C3: trigram counts

▷ length of L

▷ last two tokens in L

▷ trigram with highest count,

▷ among the trigrams starting with T

# General *n*-gram algorithm

**ngram**(L, n, Cn):                    ▷ L: tokens, Cn: *n*-gram counts

    k ← len(L)                    ▷ length of L

    T ← L[k − n + 2, ..., k]              ▷ last n − 1 tokens in L

    return argmax(Cn[T, ·])         ▷ n-gram with highest count,

                                          ▷ among the n-grams starting with T

This can replace **unigram**, **bigram** and **trigram** algorithms: just use n=1, n=2 or n=3

# Text prediction algorithm

**predict**(L, n, Cn, m): ▹ L: tokens, Cn: *n*-gram counts,

▹ m: total wanted number of words

P ← L ▹ start with words given as input

while len(P) < m: ▹ repeat until we have m words

P[len(P) + 1] ← ngram(P, n, Cn) ▹ add next word

return P ▹ output is list of words including input

# *n*-gram models comparison

| Unigram | , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , , |
|---|---|
| **Bigram** | THEN SHE WENT ON IT HAD BEEN RUNNING ABOUT IN HER HEAD ! THE GARDEN , WHO WAS NOW , FOR SOME OF THEM ! |
| **Trigram** | ALL OF A GOOD DEAL FRIGHTENED AT THE TOP OF HER SISTER , WHO WAS GENTLY BRUSHING AWAY SOME DEAD LEAVES THAT HAD FALLEN INTO A TREE A FEW MINUTES , IT WAS THE WHITE RABBIT , WHO WAS NOW ABOUT TWO FEET HIGH . |
| **4-gram** | THE FIRST THING I 'VE GOT TO DO , SO ALICE SOON BEGAN TALKING TO HERSELF . `` DINAH 'LL MISS ME VERY MUCH TO-NIGHT , I SHOULD THINK ! " ( DINAH WAS THE CAT |
| **5-gram** | AND SO IT WAS INDEED ! SHE WAS NOW ONLY TEN INCHES HIGH , AND HER FACE BRIGHTENED UP AT THE THOUGHT THAT SHE WAS NOW ABOUT TWO FEET HIGH AND WAS GOING ON SHRINKING RAPIDLY . |

# Back-off

*n*-gram models quickly become too **sparse**.

WHENEVER I WISH I

does not occur in *Alice in Wonderland*: cannot use 5-gram.

If no match is found, use a smaller *n*:

To predict the next token, **back-off** to 4-grams:

WHENEVER I WISH I COULD

| I WISH I COULD | 2 |
|---|---|
| I WISH I HAD | 2 |

# Trigram with Backoff to Bigram

**trigram-backoff-bigram**(L, C2, C3):  ▷ L: tokens,

   k ← len(L) ▷ C2: bigram counts, C3: trigram counts

   if C3[L[k − 1, k], ·]) is empty :     ▷ not found

      return argmax(C2[L[k], ·])    ▷ use bigram

   else:                  ▷ trigram found

      return argmax(C3[L[k − 1, k], ·])▷ use trigram

# Trigram with Backoff to Bigram

**trigram-backoff-bigram**(L, C2, C3):   ▷ L: tokens,

   k ← len(L) ▷ C2: bigram counts, C3: trigram counts

   if C3[L[k − 1, k], ·]) is empty :     ▷ not found

     return argmax(C2[L[k], ·])    ▷ use bigram

  else:                   ▷ trigram found

     return argmax(C3[L[k − 1, k], ·])▷ use trigram

**if/else condition**

# Trigram with Full Backoff

**trigram-backoff**(L, C):         ▷ L is a list of tokens,

     k ← len(L)         ▷ C is the list [C1, C2, C3]:

     i ← 3         ▷ unigram, bigram, trigram counts

     while C[i][L[k − i + 2, ..., k], ·] is empty:

         i ← i − 1         ▷ *i*-gram not found, try *i* − 1

     return argmax(C[i][L[k − i + 2, ..., k], ·])

# References

- Google Ngram Viewer: books.google.com/ngrams
- *Alice's Adventures in Wonderland* on Wikisource:

  en.wikisource.org/wiki/Alice's_Adventures_in_Wonderland_(1866)
- *n*-grams: en.wikipedia.org/wiki/N-gram

xkcd.com/1068