

# First-two-char input method

Orcun Ilker Doeger

## 1. Introduction

Efficient text input is essential nowadays, particularly on devices with small keyboards like smartphones. Traditional typing can also be slow and prone to mistakes, which makes finding faster and more accurate input methods even more important. One approach is the first-two-char input method, where only the first two characters of each word are needed, and the system predicts and completes the word. This project focuses on implementing this method using KyTea and Python.

We will train a model using a dataset of word pairs, each consisting of the first two characters and the full word. This report outlines the steps taken to prepare the dataset, train the model, and test its performance. By evaluating the model's accuracy, we aim to demonstrate the effectiveness of the first-two-char input method in improving typing efficiency in a statistical manner.

## 2. Method

1. **Training data:** Some training data is needed to be able to train a model. For that, we prepare the data with Python, that we then can feed into KyTea. The result is on the right:

```
1 # Makes the training data out of the whole book (book.txt)
2 # Project -> Pr/Project and so on
3 # saves it into "book_processed_train.txt"
4
5 def prepare_text_for_training(text):
6     # Split text into words
7     words = text.split()
8
9     # Clean each word and bring it into form Pr/Project
10    processed_words = []
11    for word in words:
12        clean_word = re.sub('[^a-zA-Z]', '', word) #only alphabetical
13        if len(clean_word) > 1:
14            processed_word = f"{clean_word[:2]}/{clean_word}"
15        elif len(clean_word) == 1:
16            processed_word = f"{clean_word}/{clean_word}"
17        if processed_word: # only if existend
18            processed_words.append(processed_word)
19
20    processed_words += "\n"
21    return processed_words
22
23 # Do it
24 # Read
25 with open('book.txt', 'r') as file:
26     text = file.read()
27
28 #Process
29 processed_words = prepare_text_for_training(text)
30
31 # Write
32 with open('book_processed_train.txt', 'w') as file:
33     file.write('\n'.join(processed_words))
34
```

```
to/to
nu/number
my/my
da/days
th/that
I/I
ma/may
ap/apply
my/my
he/heart
un/unto
wi/wisdom
Sa/Saturday
Ma/March
Ma/March
My/My
un/uncle
re/returned
ho/home
an/and
wi/with
hi/him
my/my
au/aunt
Bl/Bly
```

## 2. Testing data: To test the model, we prepare a 2-gram-word text out of normal text.

```
1 # Prepares the test text that i took out of the book
2 # Removes any non alphabetical characters and makes 2-grams of the words
3 # for testing
4
5 def generate2gramText(input_file, output_file):
6     # Read test input
7     with open(input_file, 'r') as file:
8         text = file.read()
9
10    words = text.split()
11
12    processed_words = []
13    for word in words:
14        clean_word = re.sub('[^a-zA-Z]', '', word)
15        if clean_word:
16            processed_words.append(clean_word[:2])
17
18    # Join the processed words with a space
19    result = ' '.join(processed_words)
20    result += "\n"
21
22    # Write
23    with open(output_file, 'w') as file:
24        file.write(result)
25
26
27 input_file = 'testing_input.txt'
28 output_file = 'testing_input_processed_test.txt'
29
30 generate2gramText(input_file, output_file)
```

Sunday 8. The Lord was pleased to call her to himself. The day before she seemed much better. But in the evening she fell worse again, and this morning, about nine, departed. When word came of this, I was not shocked at first; but in a while I began to reason on what had passed between us the evening before. She then said, "When death seemed nigh some nights since, the enemy thrust sore at me, and said, This will be thy end: But I said, Not without God's permission: He can bring me down to the belly of hell, and bring me up again, as he did \_Jonah\_ out of the whale's belly. But it is one thing to talk of death, and another

Su Th Lo wa pl to ca he to hi Th da be sh se mu be Bu in th ev sh fe wo ag an th mo  
ab ni de Wh wo ca of th I wa no sh at fi bu in a wh I be to re on wh ha pa be us th  
ev be Sh th sa Wh de se ni so ni si th en th so at me an sa Th wi be th en Bu I sa No  
wi Go pe He ca br me do to th be of he an br me up ag as he di Jo ou of th wh be Bu  
it is on th to ta of de an an to ha it br ho to us Th he hu an ch an fr ca ne an so  
se to ov he Bu sh sa Go is ab to bl an ke th wi me if he is pl to ta me An if he is

## 3. Training and Testing with KyTea

Using following commands, we can now develop a model with the training data which we can then use to predict the words from the 2-gram-word-text.

```
train-kytea -full book_processed_train.txt -nows -model train.model
```

```
orcun@MacBook-Air-von-Orcun python % train-kytea -full book_processed_train.txt -nows -model train.model
Scanning dictionaries and corpora for vocabulary
Reading corpus from book_processed_train.txt done (90324 lines)
Building dictionary index done!
Creating tagging features (tag 1) done!
Training local tag classifiers done!
Printing model to train.model done!
```

```
kytea -model train.model -out tags < testing_input_processed_test.txt > test.out
```

```
orcun@MacBook-Air-von-Orcun python % kytea -model train.model -out tags < testing_input_processed_test.txt > test.out
```

4. **Evaluation:** After having a test-result (the predictions), we can now compare the resulting text to the ground-truth text:

```
1 # This compares the test text (ground-truth) to the kytea model result (from the 2grams)
2
3 def read_clean_text(input_file):
4     with open(input_file, 'r') as file:
5         text = file.read()
6
7     cleaned_text = re.sub(r'^a-zA-Z\s', '', text).lower() # Clean text
8
9     words = cleaned_text.split() # Split to words
10    return words
11
12 def compare_texts(input_file1, input_file2):
13     # Clean both result and ground-truth texts
14     words1 = read_clean_text(input_file1)
15     words2 = read_clean_text(input_file2)
16
17     matched_words = 0
18     if len(words1) == len(words2):
19         for i in range(0, len(words1)):
20             if words1[i] == words2[i]:
21                 matched_words += 1
22
23     percentage = 100 * matched_words / len(words1)
24     return percentage
25
26
27 input_file1 = 'testing_input.txt'
28 input_file2 = 'test.out'
29
30 percentage = compare_texts(input_file1, input_file2)
31 print(f"Correctness: {percentage:.2f}%")
```

words1	words2
0 sunday	0 sunday
1 the	1 the
2 lord	2 lord
3 was	3 was
4 pleased	4 pleased
5 to	5 to
6 call	6 can
7 her	7 her
8 to	8 to
9 himself	9 his
10 the	10 the
11 day	11 day
12 before	12 be
13 she	13 she
14 seemed	14 see
15 much	15 much
16 better	16 be
17 but	17 but
18 in	18 in
19 the	19 the
20 evening	20 every

### 3. Language Resource

For the training data, a raw-text book has been used: <https://www.gutenberg.org/ebooks/73797>  
It is consisting of 90500 words.

Author	<a href="#">Wesley, John, 1703-1791</a>
Title	The works of the Rev. John Wesley, Vol. 13 (of 32)
Original Publication	Bristol: William Pine, 1771.
Credits	Richard Hulse and the Online Distributed Proofreading Team at <a href="https://www.pgdp.net">https://www.pgdp.net</a> (This file was produced from images generously made available by The Internet Archive)
Language	Englisch
LoC Class	<a href="#">BX: Philosophy, Psychology, Religion: Christianity: Churches, Church movements</a>
Subject	<a href="#">Theology -- Early works to 1800</a>
Subject	<a href="#">Methodist Church</a>
Subject	<a href="#">Theology -- History -- 18th century</a>
Category	Text
EBook-No.	73797
Release Date	08.06.2024
Copyright Status	Public domain in the USA.

## 4. Experimental Evaluation

For the testing of the model, a short part of the book consisting of 350 words has been used. This experiment yields an accuracy of 58.17%:

Correctness: 58.17%			
20	evening	20	every
21	she	21	she
22	fell	22	fear
23	worse	23	would
24	again	24	again
25	and	25	and
26	this	26	the
27	morning	27	more
28	about	28	about
29	nine	29	night
30	departed	30	dear
31	when	31	when
32	word	32	would
33	came	33	can
34	of	34	of
35	this	35	the
36	i	36	i
37	was	37	was
38	not	38	not
39	shocked	39	she
40	at	40	at
41	first	41	find
42	but	42	but
43	in	43	in
44	a	44	a
45	while	45	which

## 5. Conclusion

This implementation of the first-two-char input method using KyTea and Python achieved an accuracy of 58.17% with this specific training and test data sets. While this shows that the model can predict full words from the first two characters to some extent, there is substantial room for improvement.

The moderate accuracy proposes that predicting words based on just the first two characters is challenging, due to ambiguities and the limited context. Several factors could have influenced the model's performance, such as the size and diversity of the training dataset or the complexity of the language.