



Next Word Auto-Completion

KYU CHO

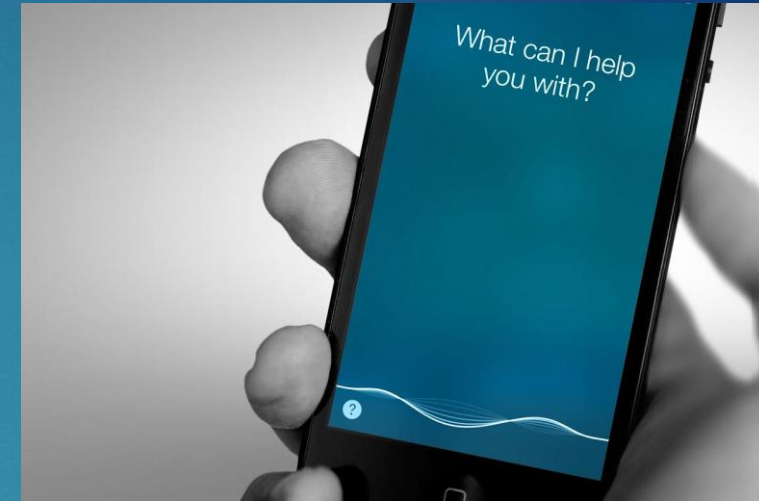
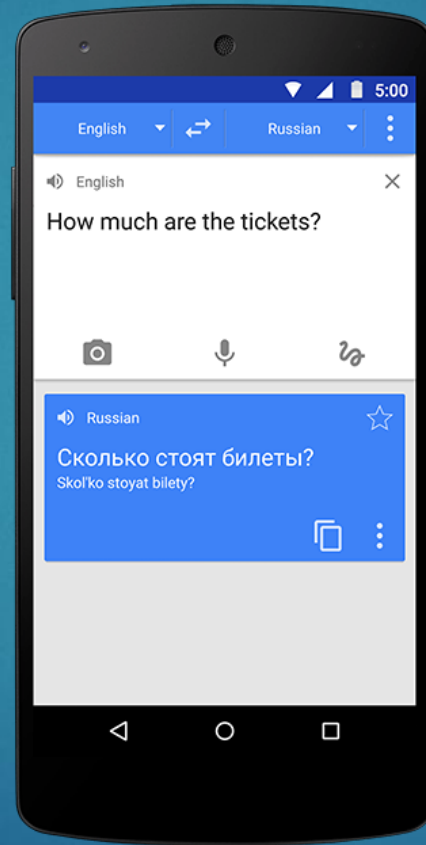
100

-



Natural Language Process (NLP) in smartphone

- ▶ Language Translation
- ▶ Sentiment Analysis
- ▶ Voice Recognition
- ▶ Question Answering
- ▶ Text Prediction
- ▶ Word Completion
- ▶ Spelling Correction
- ▶ Authorship Identification



App

<https://kyucho.shinyapps.io/nextword/>

100

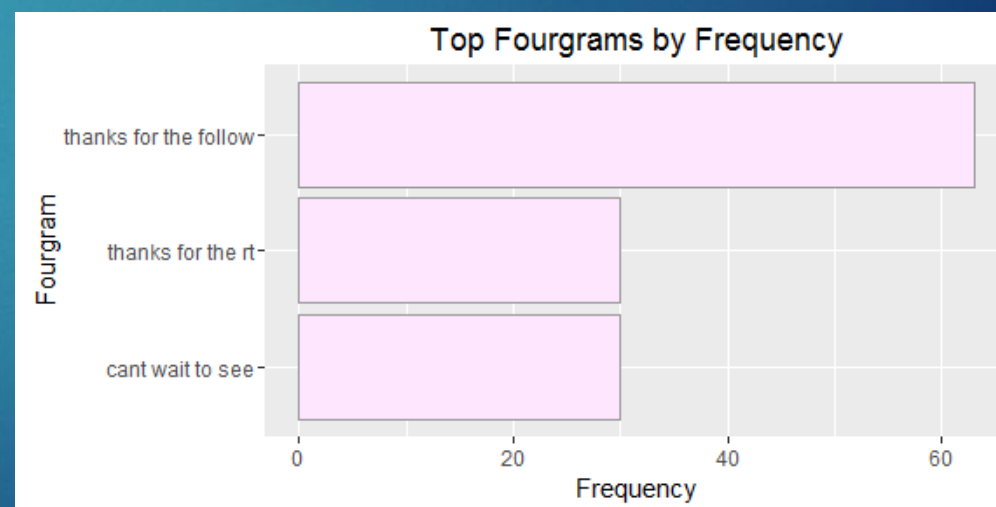
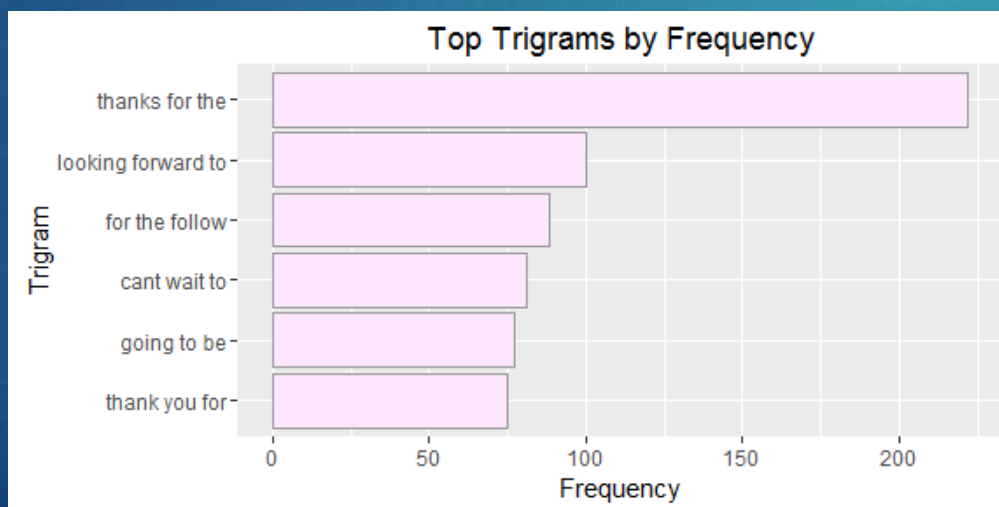
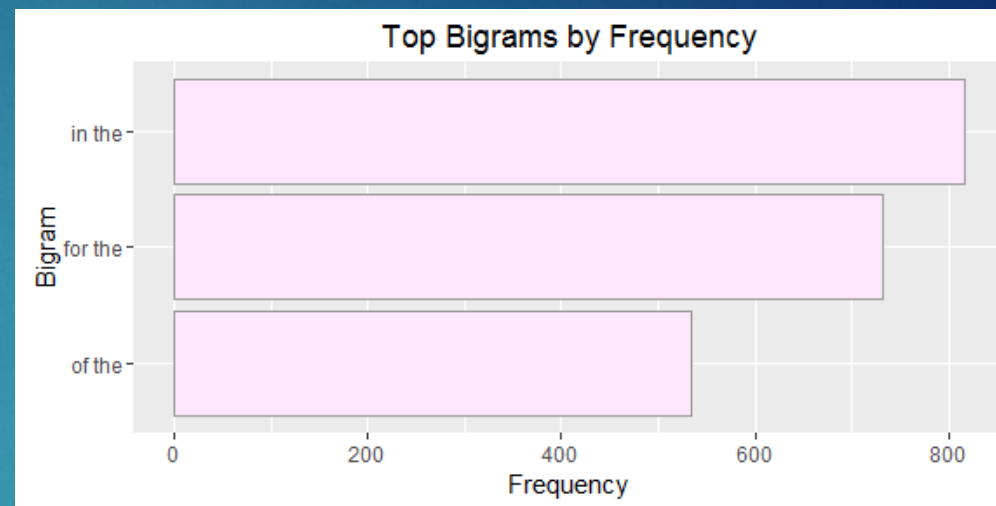
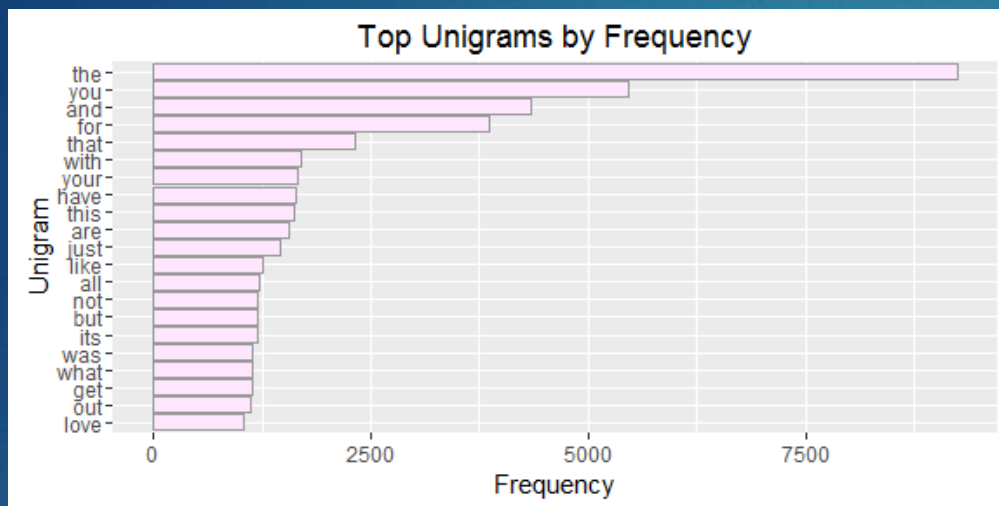
Computing the Probability of a Word Sequence

- ▶ Try: Computing the product of component conditional probabilities?
 - ▶ $P(\text{I want to eat Italian food}) = P(\text{I}) * P(\text{want} \mid \text{I}) * P(\text{to} \mid \text{I want})$
 - * $P(\text{eat} \mid \text{I want to}) * P(\text{Italian} \mid \text{I want to eat})$
 - * $P(\text{food} \mid \text{I want to eat Italian})$
- ▶ Problem: The longer the sequence, the less likely we are to find it in a training data set.
 - ▶ $P(\text{Most biologists and folklore specialists believe that in fact the mythical unicorn horns derived from the narwhal})$
- ▶ Solution: Maximum Likelihood Estimation(MLE) using N-grams

Maximum Likelihood Estimation (MLE) with N-grams

- ▶ Markov Assumption:
 - ▶ The Probability of a word depends only on the probability of a limited history.
- ▶ Generalization:
 - ▶ The probability of a word depends only on the probability of the N-previous words.
 - ▶ bi-gram, tri-grams, quad-grams,...

Plot



N-grams Examples

N-gram formula

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

- ▶ Probability with Brute Force
 - ▶ $P(\text{I want to eat Italian food}) = P(\text{I}) * P(\text{want} | \text{I}) * P(\text{to} | \text{I want})$
* $P(\text{eat} | \text{I want to}) * P(\text{Italian} | \text{I want to eat})$
* $P(\text{food} | \text{I want to eat Italian})$
- ▶ Probability of Bi-gram
 - ▶ $P(\text{I want to eat Italian food}) = P(\text{food} | \text{Italian})$
- ▶ Probability of Tri-gram
 - ▶ $P(\text{I want to eat Italian food}) = P(\text{food} | \text{Italian}) * P(\text{food} | \text{eat Italian})$
- ▶ Probability of Quad-gram
 - ▶ $P(\text{I want to eat Italian food}) =$
 $P(\text{food} | \text{Italian}) * P(\text{food} | \text{eat Italian}) * P(\text{food} | \text{to eat Italian})$

Training and Testing

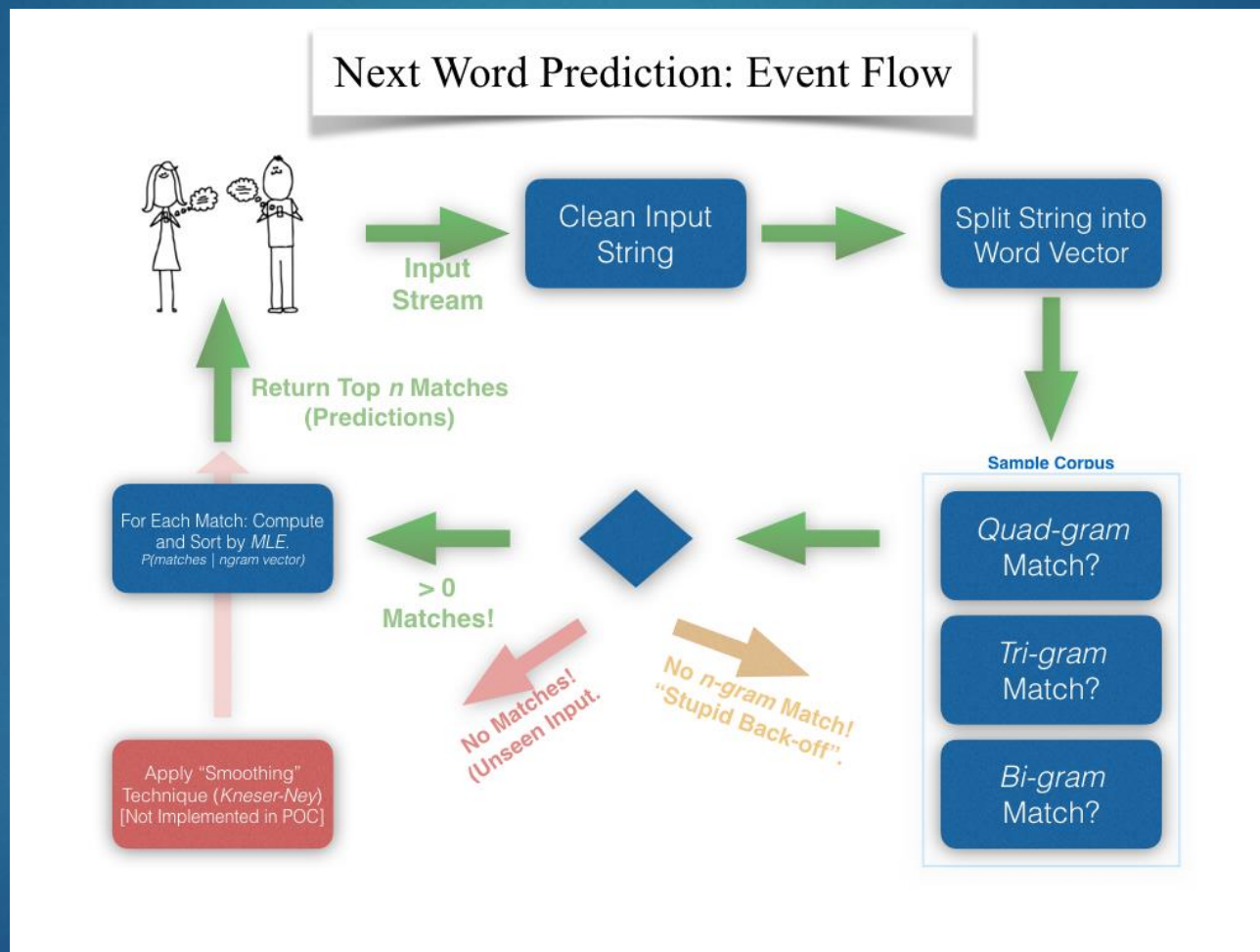
- ▶ Overly narrow corpus: probabilities don't generalize
 - ▶ Ex) brute force
- ▶ Overly general corpus: probabilities don't reflect task or domain
 - ▶ Ex) Uni-gram, bi-gram

Backoff

- ▶ The higher N-gram needs more data to train.
Thus backoff models is needed for missing value.
- ▶ **Stupid-backoff** – simple yet powerful
If the quad-gram has no match in the data set, move to tri-gram
If the tri-gram has no match in the data set, move to bi-gram. Etc

“want to eat” in quad-gram data set? No then
“to eat” in tri-gram data set? No then
“eat” in bi-gram data set? yes then
what’s the predicted word after “eat”?

Flow Chart

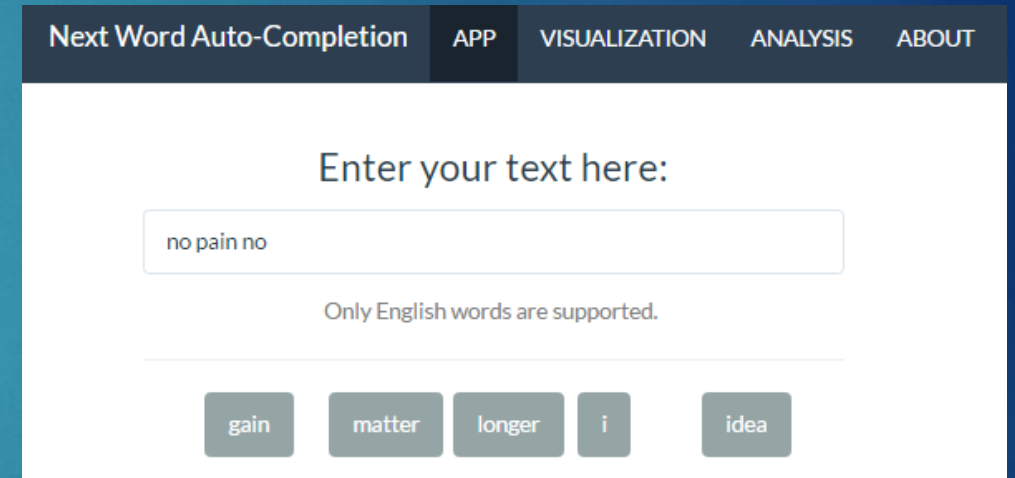


Table

Prediciton Table			
Show 10 entries		Search: <input type="text"/>	
input	output	smle	counts.y
the most important	thing	0.2008	148
important	to	0.1692	1331
most important	thing	0.1568	164
the most important	part	0.0488	36
the most important	things	0.0434	32
most important	to	0.0402	42
most important	part	0.0402	42
important	than	0.0379	298
important	for	0.0371	292
most important	things	0.0354	37
<input type="text" value="input"/>	<input type="text" value="output"/>	<input type="text" value="smle"/>	<input type="text" value="counts.y"/>
Showing 1 to 10 of 80 entries			
<div>Previous12345...8Next</div>			

What have I done for past few months?

- ▶ Collect data from twitter API (0.5 GB Data)
- ▶ Clean data
 - ▶ remove bad words and symbols,
 - ▶ won't -> will not, etc
- ▶ Create N-gram up to 4
- ▶ Calculate Maximum Likelihood Estimation
- ▶ Data Explanatory Analysis
- ▶ Build the Model based on the N-gram data set with MLE
- ▶ Develop Front-end and Back-end for UI



The screenshot shows a web application titled "Next Word Auto-Completion". The navigation bar includes links for "APP", "VISUALIZATION", "ANALYSIS", and "ABOUT". The main content area has a heading "Enter your text here:" followed by a text input field containing "no pain no". Below the input field, a message states "Only English words are supported." At the bottom, there are five buttons displaying suggested words: "gain", "matter", "longer", "i", and "idea".

THANK YOU

- ▶ App

- ▶ <https://kyucho.shinyapps.io/nextword/>

- ▶ Source

- ▶ <https://github.com/jamin567/DataScienceCapston>

- ▶ Inquiry

- ▶ chok20734@gmail.com