

Natural Language Processing Lab

Week3: UNIX Lab

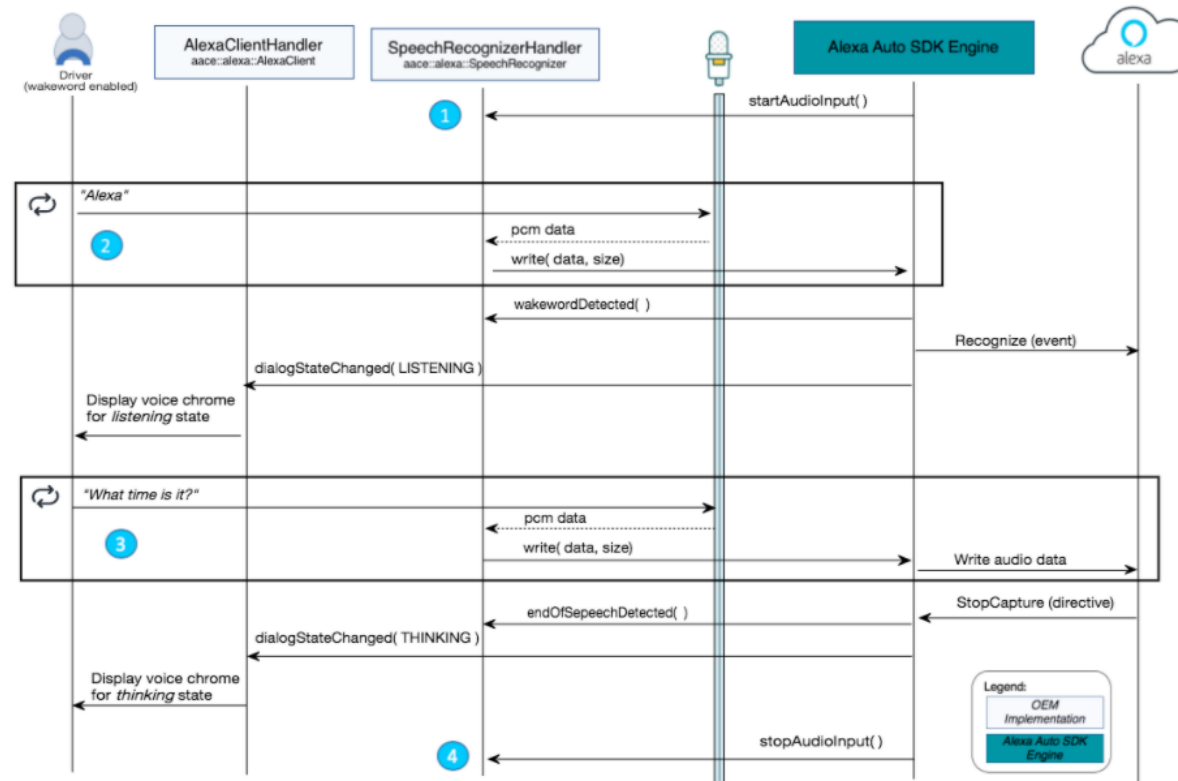
Praveen Joshi

05/10/2020

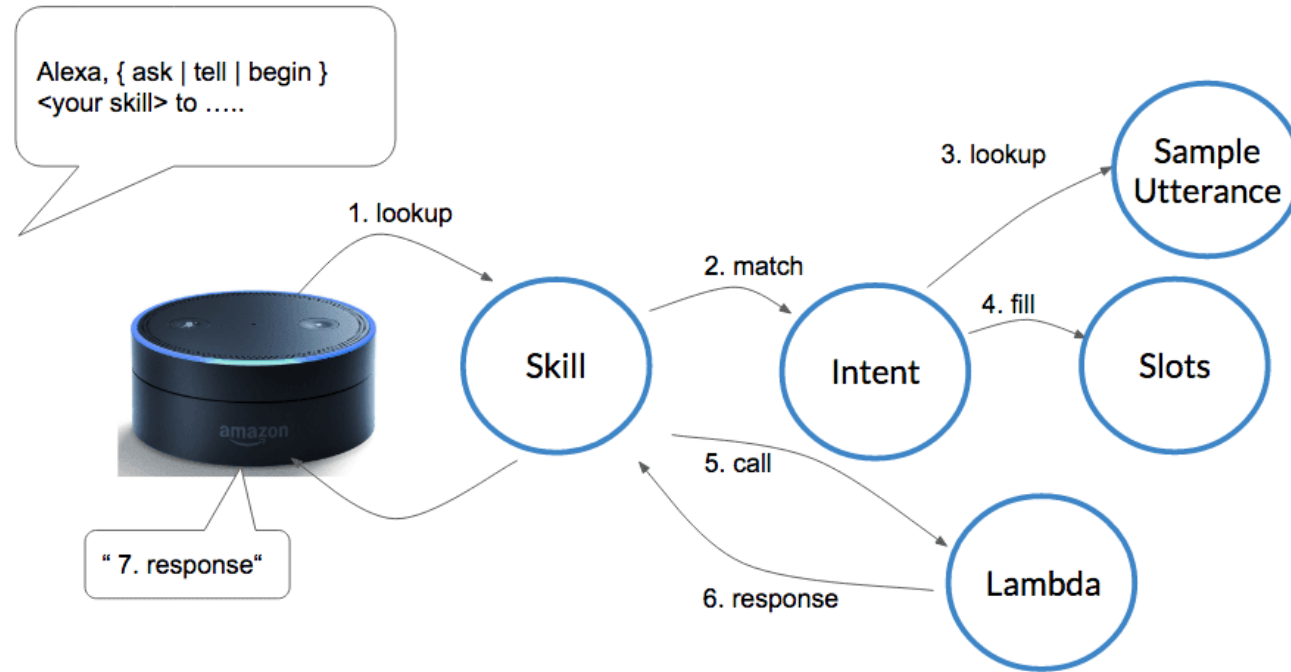
Alexa Custom Skill - Reference Architecture



<https://www.krminc.com/alexa-app-development/>



https://alexa.github.io/alexa-auto-sdk/SEQUENCE_DIAGRAMS.html



<https://www.dipockdas.com/2016/11/create-alexa-skill-for-your-amazon.html>

- Why?
 - Matching/Finding
 - Doing something with matched text
 - Validation of data
 - Case insensitive matching
 - Parsing data (ex: html)
 - Converting data into diff for etc.

■ Entities to extract:

- Board of Directors
- Price at NAV
- Fund Assets (millions)
- Expense Ratio
- Minimum Initial Investment
- Fund Name, Class and series extraction
- Annual Management Fees
- Expense Ratio

ATLANTE FUNDS PLC

DIRECTORS

Peter Blessing
Dermot Butler
Stuart Anthony Williams
Matteo Riginello

SPONSOR

Albemarle Asset Management Limited
28-29 Dover Street
London W1S 4NA
United Kingdom

REGISTERED OFFICE

70 Sir John Rogerson's Quay
Dublin 2
Ireland

INVESTMENT MANAGER

Albemarle Asset Management Limited
28-29 Dover Street
London W1S 4NA
United Kingdom

SECRETARY

Matsack Trust Limited
70 Sir John Rogerson's Quay
Dublin 2

ADMINISTRATOR

Bank of Ireland Securities Services Limited
New Century House
Mayor Street Lower
IFSC
Dublin 1
Ireland

https://github.com/praveenjoshi01/Hedge_Fund_Information_Extraction

- The process we just went through was based on fixing two kinds of errors
 - Matching strings that we should not have matched (there, then, other)
 - False positives (Type I)
 - Not matching things that we should have matched (The)
 - False negatives (Type II)

Slide taken from Dr. Haithem Afli's Lecture on NLP

Confusion Matrix



		Actual Value (as confirmed by experiment)	
		positives	negatives
Predicted Value (predicted by the test)	positives	TP True Positive	FP False Positive
	negatives	FN False Negative	TN True Negative

<https://alearningaday.blog/2016/09/14/confusion-matrix/>

Confusion Matrix



- Set of words: {the, The, There, then, other, there}
- Ground Truth: {the, The}
- **Regex**(Set of words): {the, there, then, other}

		Actual Value (as confirmed by experiment)	
		positives	negatives
Predicted Value (predicted by the test)	positives	TP True Positive	FP False Positive
	negatives	FN False Negative	TN True Negative

<https://alearningaday.blog/2016/09/14/confusion-matrix/>

	Actual Value	
	Positives	Negatives
Predicted Value		
Positives	the	there, then, other
Negatives	The	There

Confusion Matrix – Regex Expressions



- Matching strings that we should not have matched (**there**, **then**, **other**)
 - False positives (Type I)
- Not matching things that we should have matched (**The**)
 - False negatives (Type II)

		Actual Value (as confirmed by experiment)	
		positives	negatives
Predicted Value (predicted by the test)	positives	TP True Positive	FP False Positive
	negatives	FN False Negative	TN True Negative

<https://alearningaday.blog/2016/09/14/confusion-matrix/>

	Actual Value		
Predicted Value		Positives	Negatives
	Positives	the	there, then, other
	Negatives	The	There

Slide taken from Dr. Haithem Afli's Lecture on NLP

Confusion Matrix

		CONDITION determined by "Gold Standard"			
TOTAL POPULATION		CONDITION POS	CONDITION NEG	PREVALENCE $\frac{\text{CONDITION POS}}{\text{TOTAL POPULATION}}$	
TEST OUT- COME	TEST POS	True Pos TP	Type I Error False Pos FP	Precision Pos Predictive Value $\text{PPV} = \frac{\text{TP}}{\text{TEST P}}$	False Discovery Rate $\text{FDR} = \frac{\text{FP}}{\text{TEST P}}$
	TEST NEG	Type II Error False Neg FN	True Neg TN	False Omission Rate $\text{FOR} = \frac{\text{FN}}{\text{TEST N}}$	Neg Predictive Value $\text{NPV} = \frac{\text{TN}}{\text{TEST N}}$
ACCURACY ACC $\text{ACC} = \frac{\text{TP} + \text{TN}}{\text{TOT POP}}$		Sensitivity (SN), Recall Total Pos Rate TPR $\text{TPR} = \frac{\text{TP}}{\text{CONDITION POS}}$	Fall-Out False Pos Rate FPR $\text{FPR} = \frac{\text{FP}}{\text{CONDITION NEG}}$	Pos Likelihood Ratio LR + $\text{LR} + = \frac{\text{TPR}}{\text{FPR}}$	Diagnostic Odds Ratio DOR $\text{DOR} = \frac{\text{LR} +}{\text{LR} -}$
		Miss Rate False Neg Rate FNR $\text{FNR} = \frac{\text{FN}}{\text{CONDITION POS}}$	Specificity (SPC) True Neg Rate TNR $\text{TNR} = \frac{\text{TN}}{\text{CONDITION NEG}}$	Neg Likelihood Ratio LR - $\text{LR} - = \frac{\text{TNR}}{\text{FNR}}$	

<https://www.unite.ai/what-is-a-confusion-matrix/>

Confusion Matrix



- Set of words: {the, The, There, then, other, there}
- Ground Truth: {the, The}
- **Regex**(Set of words): {The, There}
- Identify Quadrant?
 - Q1
 - Q2
 - Q3
 - Q4
 - Type 1 error
 - Type 2 error

	Predicted Value		
Actual Value		Positives	Negatives
	Positives	Q1	Q2
	Negatives	Q3	Q4

Confusion Matrix



- Set of words: {the, The, There, then, other, there}
- Ground Truth: {the, The}
- **Regex**(Set of words): {The, There}
- Identify Quadrant?
 - Q1
 - Q2
 - Q3
 - Q4
 - Type 1 error
 - Type 2 error

	Predicted Value		
Actual Value		Positives	Negatives
	Positives	TP	FN/T2
	Negatives	FP/T1	TN

Confusion Matrix



- Set of words: {the, The, There, then, other, there}
- Ground Truth: {the, The}
- **Regex**(Set of words): {The, There}
- Identify Quadrant?
 - the
 - The
 - There
 - Then
 - other
 - there

	Predicted Value		
		Positives	Negatives
	Actual Value		
		Positives	Negatives
		Q1	Q2
		Q3	Q4

Confusion Matrix



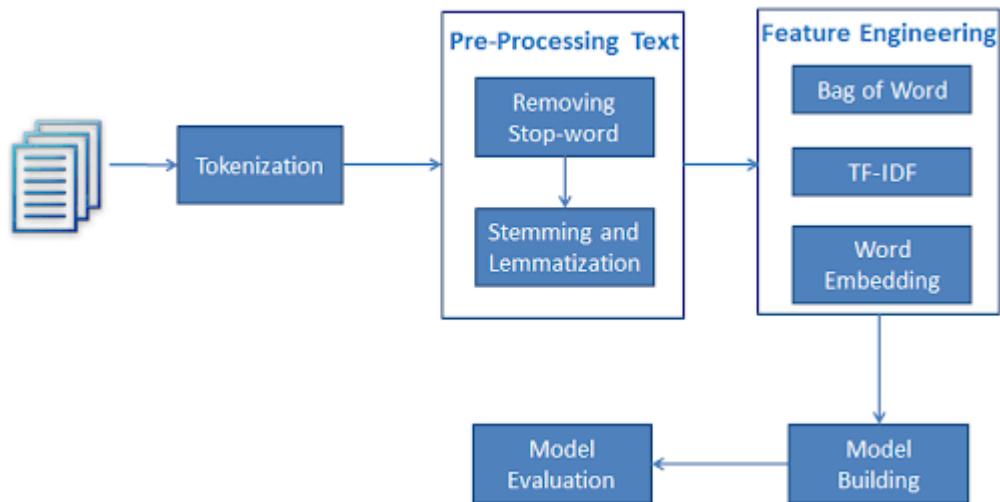
- Set of words: {the, The, There, then, other, there}
- Ground Truth: {the, The}
- **Regex**(Set of words): {The, There}
- Identify Quadrant?
 - the
 - The
 - There
 - Then
 - other
 - there

	Actual Value		
		Positives	Negatives
	Positives	The	There
	Negatives	the	then, other, there

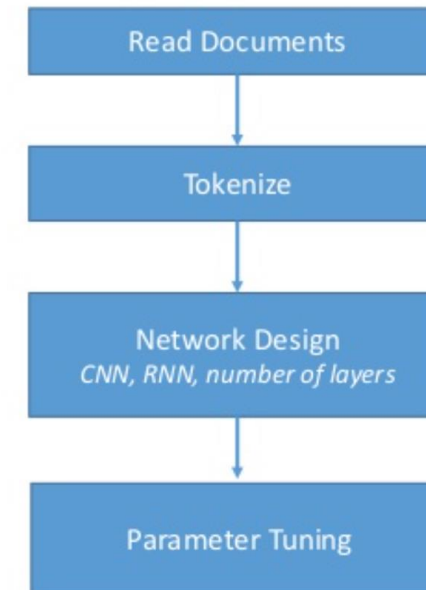
- In NLP we are always dealing with these kinds of errors.
- Reducing the error rate for an application often involves two antagonistic efforts:
 - Increasing accuracy or precision (minimizing false positives)
 - Increasing coverage or recall (minimizing false negatives).

Slide taken from Dr. Haithem Afli's Lecture on NLP

■ Steps of Text Classification



Traditional Approach



Deep Learning Approach

<https://www.slideshare.net/SomnathBanerjee17/classifying-text-with-cnn>
<https://www.datacamp.com/community/tutorials/text-analytics-beginners-nltk>

- **Cygwin**

Thank you

Praveen Joshi

05/10/2020