
SWAHILI CHATBOT IMPLEMENTATION

PROJECT REPORT

Amittai Siavava
Dartmouth College
amittai.j.wekesa.24@dartmouth.edu

January 12, 2023

ABSTRACT

This project simulates an eliza bot using regular expressions. It interprets user input in Swahili language and responds accordingly.

1 Supported Statements

All statements were adapted from the original Spanish statements in the assignment.

1.1 Initial Greeting

Prompt		Response	
<i>Jina langu ni Amittai.</i>	My name is Amittai.	Shikamoo, Amittai.	Hello, Amittai. How have you been?
<i>Naitwa Amittai.</i>	I am called Amittai.	Umeshindaje?	

Table 1: Greetings and Responses.

1.2 State of Mind

Prompt		Response	
<i>Nimefurahi.</i>	I am happy.	<i>Mbona umefurahi?</i>	Why are you happy?
<i>Mimi nimefurahi.</i>	*I* am happy.		
<i>Nimefurahi sana.</i>	I am very happy.		
<i>Sijafurahi.</i>	I am not happy.	<i>Mbona hujafurahi?</i>	Why aren't you happy?
<i>Sijafurahi sana.</i>	I am not very happy.		
<i>Mimi sijafurahi.</i>	*I* am not happy.		
<i>Nimehuzunika.</i>	I am sad.	<i>Mbona umehuzunika?</i>	Why are you sad?
<i>Mi nimehuzunika kidogo.</i>	I'm a little bit sad.		
<i>Nimehuzunika kidogo.</i>	I am a little bit sad.		
<i>Nimehuzunika kiasi.</i>	I am quite sad.		
<i>Nimehuzunika kidogo.</i>	I am sort of sad.		

Table 2: State-of-Mind Statements.

1.3 Characteristics of a Person

Prompt		Response	
<i>Mimi ni mchafu.</i>	I am untidy.	<i>Mbona wewe ni mchafu?</i>	Why are you untidy?
<i>Ni kwamba mimi ni mchafu sana.</i>	It's just that I am fairly untidy.	<i>Mbona wewe ni mchafu sana?</i>	Why are you fairly untidy?
<i>Mimi kawaida ni mtu mchafu.</i>	I am usually dirty	<i>Mbona wewe kawaida ni mtu mchafu?</i>	Why are you usually dirty?
<i>Mimi kawaida si mtu safi.</i>	I am not happy.	<i>Mbona wewe kawaida si mtu safi?</i>	Why are you usually not a clean person?
<i>Mimi ni mwanafunzi wa Dartmouth.</i>	I am a student at Dartmouth.	<i>Mbona wewe ni mwanafunzi wa Dartmouth?</i>	Why are you a student at Dartmouth?

Table 3: Person Descriptions and Responses.

1.4 About Family

Note: I expanded the regex implementation here to also include *kakako* (“your brother”) and *dadako* (“your sister”).

Prompt		Response	
<i>Mamangu anapenda waridi.</i>	My mom likes roses.	<i>Niambie mengine kumhusu mamako.</i>	Tell me more about your mother.
<i>Babangu hapendi kakati.</i>	My dad does not like cacti.	<i>Niambie mengine kumhusu babako.</i>	Tell me more about your father.
<i>Kwa sababu mamangu yu hodari na mimea.</i>	Because my mother is good with plants.	<i>Niambie mengine kumhusu mamako.</i>	Tell me more about your mother.
<i>Kwa sababu babangu anapendelea miti.</i>	Because my father prefers trees.	<i>Niambie mengine kumhusu babako.</i>	Tell me more about your father.

Table 4: Family Descriptions.

1.5 Modal Verbs

Prompt		Response	
<i>Nataka kujifunza kuandika msimbo.</i>	I want to learn how to program.	<i>Mbona unataka kujifunza kuandika msimbo?</i>	Why do you want to learn how to program?
<i>Sitaki kujifunza kuandika msimbo.</i>	I don't want to learn how to program.	<i>Mbona hutaki kujifunza kuandika msimbo?</i>	Why don't you want to learn how to program?
<i>Nataka kuenda safari.</i>	I want to go on a trip.	<i>Mbona unataka kuenda safari?</i>	Why do you want to go on a trip?
<i>Sitaki kuenda safari.</i>	I don't want to go on a trip.	<i>Mbona hutaki kuenda safari?</i>	Why don't you want to go on a trip?
<i>Lazima niende kwenye sinema.</i>	I must go to the movies.	<i>Mbona lazima uende kwenye sinema?</i>	Why must you go to the movies?
<i>Lazima nihamie kwenye nyumba tofauti.</i>	I must move to a different house.	<i>Mbona lazima uhamie kwenye nyumba tofauti?</i>	Why must move to a different house?
<i>Naweza kutembea kwenye theluji.</i>	I can walk on the snow.	<i>Mbona unaweza kutembea kwenye theluji?</i>	Why can you walk on the snow?
<i>Naweza kula vidakuzi kumi.</i>	I can eat ten cookies.	<i>Mbona unaweza kula vidakuzi kumi?</i>	Why can you eat ten cookies?
<i>Rafiki yangu anazema naweza kula vidakuzi kumi.</i>	My friend says I can eat ten cookies.	<i>Mbona unaweza kula vidakuzi kumi?</i>	Why can you eat ten cookies?

Table 5: Modal Verbs and Responses.

1.6 Thoughts and Hopes

Prompt		Response	
<i>Nadhani rafiki zangu walifanya hivo pia.</i>	I think my friends did that too.	<i>Mbona unadhani hivo?</i>	Why do you think so?
<i>Natumai kuhitimu mwaka ujao.</i>	I hope to graduate next year.	<i>Mbona unatumai hivo?</i>	Why do you hope so?

Table 6: Thoughts, Hopes, and Responses.

1.7 Asking for Examples

Prompt		Response	
<i>Yeye husema hivo tu.</i>	He always says that.	<i>Unaweza nipa mfano?</i>	Can you give me an example?
<i>Kakangu hula vidakuzi.</i>	My brother always eats cookies.	<i>Unaweza nipa mfano?</i>	Can you give me an example?

Table 7: Prompting for Examples.

1.8 Handling Insults

Prompt		Response	
<i>Hapana, mshenzi!</i>	No, stupid!	<i>Ha! Matusi hapana! Tulia na uniambie mengine.</i>	Hey! No insults! Relax and tell me more.
<i>Wewe ni mjinga.</i>	You are an idiot.		

Table 8: Handling Insults.

1.9 All Other Statements

Prompt		Response	
<i>Ninapenda kusoma.</i>	I like reading.	<i>Niambie mengine...</i>	Tell me more...
<i>Mimi natazama runinga kwa zaidi.</i>	I watch too much television.		

Table 9: Other Statements.

2 Discussion

This program is obviously very limited. After making the program, please answer the following question in the form of a comment at the end of the code: What changes would you make to the program so that it can carry out a more human-like conversation? Please make three suggestions for how to improve the program and find a section of the Jurafsky and Martin textbook that can help you implement those changes. Mention those sections explicitly. The answer should be at least 200 words long.

- 1. Use a Better Model:** A main issue is that regular expressions do not have the capacity to fully understand nuances in language, such as emotionality, sarcasm, and other forms of non-literal speech. Regular expressions only account for the phrases they have been explicitly programmed to handle. A better model such as a neural network or especially a transformer would be able to understand the nuances in language better and generate more human-like responses. Transformer models and LSTM neural networks are particularly useful for language understanding and generation since they can understand the context of a sentence and the context of each word in a sentence. **Jurafsky and Martin** discuss such neural network models in sections 9 and 10.

Aside: A transformer model would also be more computationally intensive. If there are significant computational constraints, that could be a reason to stick with regular expressions.

- 2. Increase the range of responses:** The program, as is, is limited to a small set of responses to an equally small set of prompts. Increasing the repertoire of responses would make for a more human-like conversation. I think techniques such as *normalization* using *stemming* and *lemmatization* could be useful *before* matching regular expressions – so that related words (for example, “go” and “went”) are matched by the same regular expression. **Jurafsky and Martin** discuss how to do such normalization in Section 2.4.
- 3. Adding Variance:** Another shortfall, currently, is that the program *always* generates the same response for a given prompt. This is “OK” – but not human-like. Adding some noise in the decision-making process would be useful, such as by using a probabilistic model to choose a response *after matching a regular expression*. One way to achieve this could be to use a *Markov Chain*. **Jurafsky and Martin** do not discuss Markov Chains explicitly, but they discuss *Hidden Markov Models* in Section 8.4.