



THE UNIVERSITY OF
MELBOURNE

Not just another game: Hangman

Importance of Language Models

By:

Daniel Gill

Kartik Kishore

Kaushik Ramesh



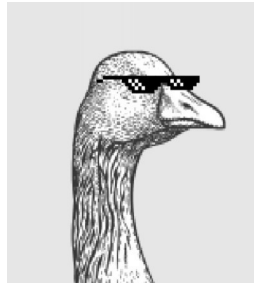


About Us



Daniel Gil

- Pursuing Masters of data science since 2018
- Extensive business process streamlining & BI experience across banking, retail, manufacturing & healthcare domains.
- Expert in Data Analytics, Business Intelligence & Natural Language Processing.



Kartik Kishore

- Ms. Data Science, University of Melbourne, Batch of 2019
- Have worked as an Identity Admin at *Cognizant* in Cyber security projects, which has been the driving force behind my interest in data analytics
- Hands-on experience in Python, R & shell. Automation & computer networks catch my fancy
- Likes pineapple on Pizza



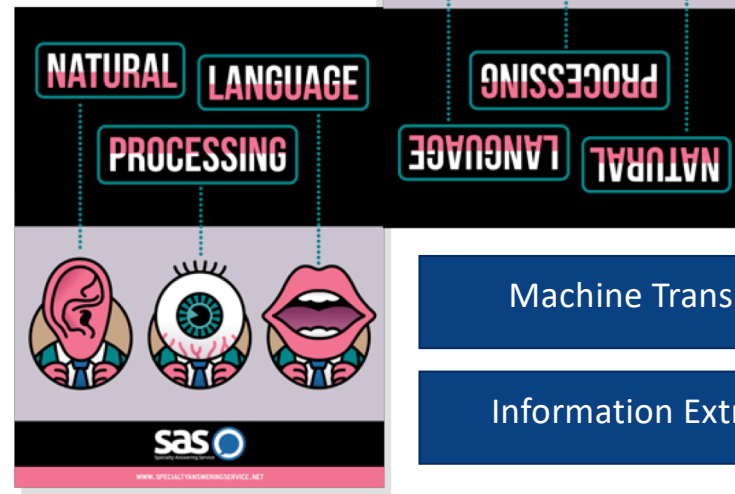
Kaushik Ramesh

- Pursuing Masters of data science since 2018
- Developed automation tools to eliminate manual interference during my tenure at Infosys limited in Bangalore, India
- Hold expertise in python, Java, R and Unix shell scripting

Introduction

Spelling Correction

Text To Speech



Machine Translation

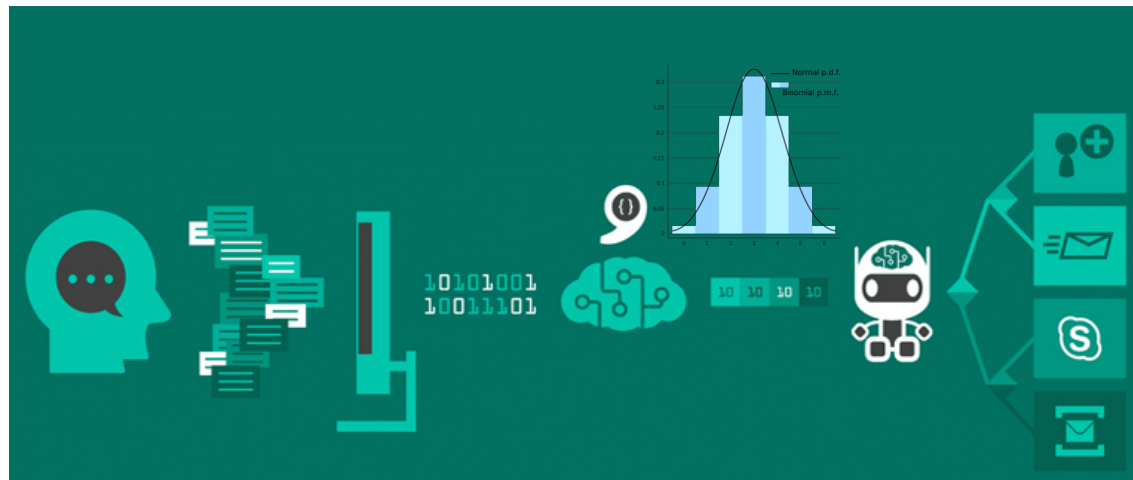
Information Extraction

Image Source: <https://www.specialtyansweringservice.net>



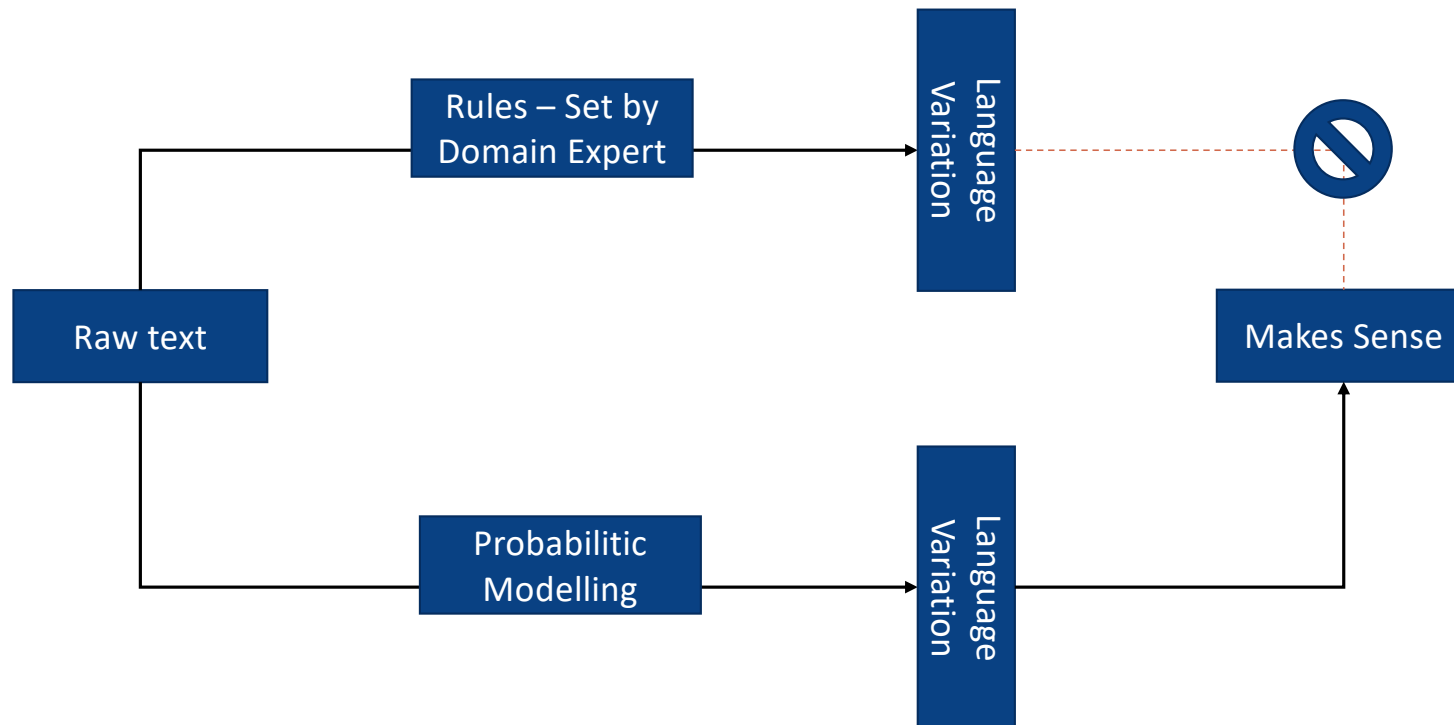
Introduction

- These are the techniques and methods involved in manipulation of any language used by us humans.
- Natural Language Toolkit (NLTK) is a package of algorithms for symbolic and statistical natural language processing for English written in the Python programming language.
- At the core of these algorithms is hardcore statistics.





Natural Language Processing





Basic Language Models

What is a language model?

A model created by following a probability distribution over strings of text

- How likely is a given string (observation) in a given corpus.
- It creates involve making an n^{th} **order Markov assumption** and **estimating n-gram probabilities** via counting and subsequent smoothing.

MARKOV ASSUMPTION: Calculations on state 'k' are influenced by output of calculation from state 'k-1'

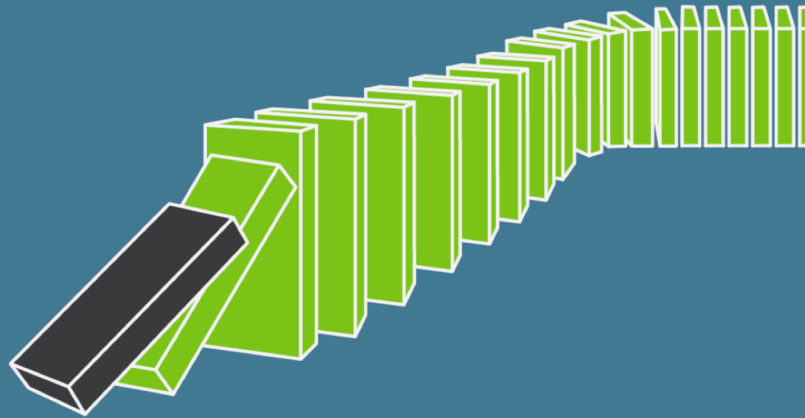
ESTIMATING N-GRAM PROBABILITIES: probability $P(w_1, w_2, \dots, w_n)$ is a product of word probabilities based on a history of preceding words, whereby the history is limited to m words

$$p(w_n | w_1, w_2, \dots, w_{n-1}) \approx p(w_n | w_{n-m}, \dots, w_{n-2}, w_{n-1})$$



Contemplations

BAD DATA = BAD EVERYTHING



Garbage in, Garbage out



Above & Beyond

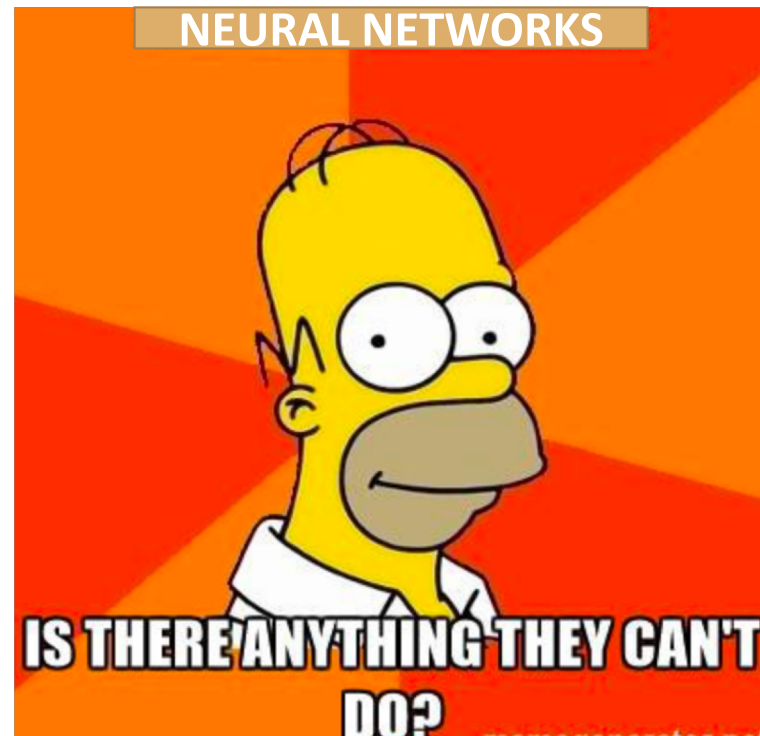


Image-source: <https://blog.goodaudience.com/introduction-to-convolution-neural-networks-ef4d03a5ca69>



State of the art

Continuous-space LM, also known as neural language model (NLM) have replaced traditional n-gram based models and have recently proved to give state of the art performance.

Varieties:

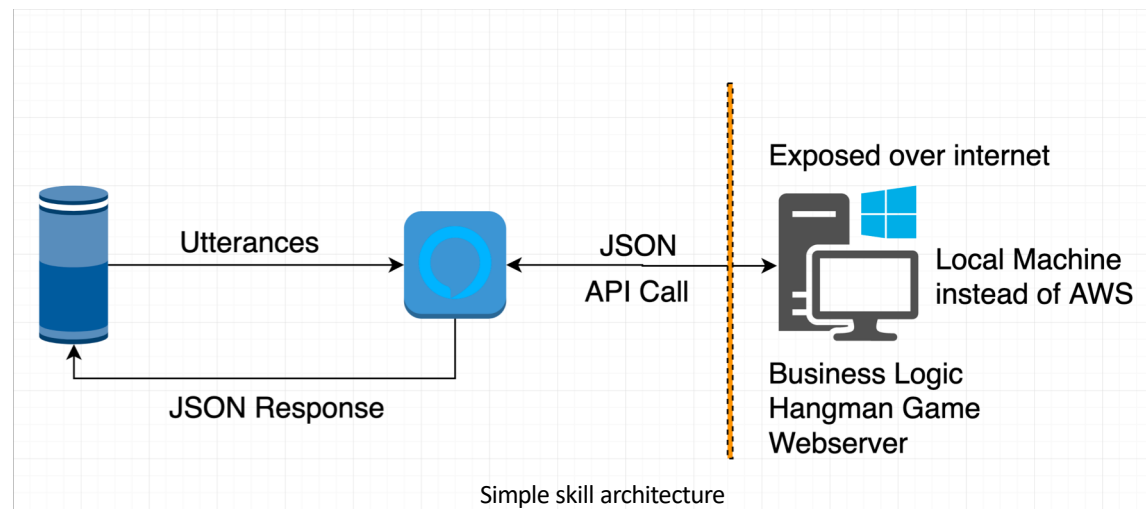
- **Feed-forward neural network based LM**
 - Which was proposed to tackle the problems of data sparsity
- **Recurrent neural network based LM**
 - Which was proposed to address the problem of limited context.



DEMO



Image Source: developer.amazon.com/console/alexa/ask





Applications in Research

- Politics researcher
- Journalism
- Human behavior analysis



Q&A

Any questions?



THE UNIVERSITY OF
MELBOURNE

Thank You