

# Predictive Input Methods

- Anish Patil

RedHat i18n team



# Agenda:

---

1. What/Why?
2. Language models
3. Projects
4. Demonstration







# Need

---

- 1.21 Billion population (India)
  - 74% literate (read & write any language)
  - Still only 5-6% understand English
  - 51% youth in 1.21 Billion
- Diversity in India
  - 22 Officially recognized languages
  - 9 Major scripts



# Rest of the world

---

- List of extinct language's
  - [http://en.wikipedia.org/wiki/List\\_of\\_extinct\\_languages\\_of\\_Europe](http://en.wikipedia.org/wiki/List_of_extinct_languages_of_Europe)
  - <http://www.unesco.org/culture/languages-atlas/en/atlasmap.html>



# Problems with Natural Languages

---

- Ambiguous
- Exceptions!
- Humans- Ignore!





# How to predict next word?

---

- Statistical Techniques
- Probability



# Probabilities

- Bayes' theorem
- $p( A | B ) = p( A ) * p( B | A ) / p( B )$
- Max has two coins in his pocket, a fair coin (head on one side and tail on the other side) and a two-headed coin. He picks one at random from his pocket, tosses it and obtains head. What is the probability that he flipped the fair coin?







$$P(\text{head}|\text{fair coin}) = \frac{1}{2}$$

$$P(\text{head}|\text{unfair coin}) = 1$$

$$P(\text{fair coin}) = \frac{1}{2}$$

$$P(\text{unfair coin}) = \frac{1}{2}$$

$$\begin{aligned} P(\text{head}) &= P(\text{head}|\text{fair coin})P(\text{fair coin}) + P(\text{head}|\text{unfair coin})P(\text{unfair coin}) \\ &= \frac{1}{2} \cdot \frac{1}{2} + 1 \cdot \frac{1}{2} \\ &= \frac{1}{4} + \frac{2}{4} = \frac{3}{4} \end{aligned}$$

$$\begin{aligned} P(\text{fair coin}|\text{head}) &= \frac{P(\text{head}|\text{fair coin})P(\text{fair coin})}{P(\text{head})} \\ &= \frac{\frac{1}{2} \cdot \frac{1}{2}}{\frac{3}{4}} \\ &= \frac{1}{4} \cdot \frac{4}{3} = \frac{1}{3} \end{aligned}$$



# Language Model

- Lot of words in one language but what is the probability that one word follow another word?
- Simple model:- number of occurrence of word/Number of words in the language



# Language Model

- A language model consists of a finite set  $V$ , and a function  $p(x_1, x_2, \dots, x_n)$  such that:

1. For any  $x_1 \dots x_n \in V^+$ ,  $p(x_1, x_2, \dots, x_n) \geq 0$

2. In addition,

$$\sum_{x_1 \dots x_n \in V^+} p(x_1, x_2, \dots, x_n) = 1$$

$$x_1 \dots x_n \in V^+$$

Hence  $p(x_1, x_2, \dots, x_n)$  is a probability distribution over the sentences in  $V^+$





# Markov Models

---

- The probability of a word depends only on the probability of a limited history
- The probability of a word depends only on the probability of the  $n$  previous words
- Unigrams, Bigrams, Trigrams...



# Markov Models cont..

- English words  $W = w_1, w_2, w_3, \dots, w_n$
- $p(w_1, w_2, w_3, \dots, w_n) = p(w_1) p(w_2|w_1) p(w_3|w_1, w_2) \dots p(w_n|w_1, w_2, \dots, w_{n-1})$
- Bigram model:-  $p(w_1, w_2, w_3, \dots, w_n) = p(w_1) p(w_2|w_1) p(w_3|w_2) \dots p(w_n|w_{n-1})$
- $w_{n-1}$  is called the history
- For example, the dog barks STOP
- $p(\text{the dog barks STOP}) = p(\text{the}|*, *) \times p(\text{dog}|*, \text{the}) \times p(\text{barks}|\text{the}, \text{dog}) \times p(\text{STOP}|\text{dog}, \text{barks})$



# Markov Models cont..

- maximum likelihood estimation
- $p(w_2|w_1) = \text{count}(w_1, w_2) / \text{count}(w_1)$





# Example

- Training Set:

START ASIA IS AWESOME STOP

START GNOME IS AWESOME STOP

START GNOME ASIA IS AWESOME STOP

- $v = \{\text{START, GNOME, .Asia, IS, AWESOME, END}\}$
- Unigram Model:-  $p(\text{ASIA}) = 2/10 = 0.2$



# Example cont..

- Trigram Model:-

$$P(\text{GNOME/START,START}) = P(2/3)$$

- $P(\text{START GNOME ASIA IS AWESOME STOP}) = P(\text{GNOME/START,START}) * P(\text{ASIA/GNOME,START}) * P(\text{IS/GNOME,ASIA}) * P(\text{AWESOME/ASIA,IS}) * P(\text{STOP/IS,AWESOME})$   
 $= (2/3) * (2/3) * (2/1) * (3/1) * (3/2)$



# Training

---

- Data??
- Testing?





# How to evaluate a L.M?

- Perplexity
- $H(W) = 1/n \log p(W)$
- Lower is the perplexity higher is accuracy of your language model



# Unseen Sentences

---

- GNOME SHELL is AWESOME
- Smoothing?
- Discounting?



# Smoothing

- Zero probabilities of unigram costs zero probabilities of entire sentence
- For unigrams add 1 to every word and adjust the count and divide it by size vocabulary to normalize it
- Original  $P(w) = c / N$
- New  $P(w) = (c+1) / (V+N)$





# Linear Interpolation

- $q_M L(w|u, v) = c(w, u, v) / c(u, v)$
- $q_M L(w|v) = c(v, w) / c(v)$
- $q_M L(w) = c(w) / c()$
- $q(w|u, v) = \lambda_1 \times q_M L(w|u, v) + \lambda_2 \times q_M L(w|v) + \lambda_3 \times q_M L(w)$
- $\lambda_1 \geq 0, \lambda_2 \geq 0, \lambda_3 \geq 0$  and  $\lambda_1 + \lambda_2 + \lambda_3 = 1$



# Libyokan and libyokan- data

---

- <https://gitorious.org/libyokan>
- <https://gitorious.org/yokan-data-mr-in>



# ibus-typing-booster

- <https://fedorahosted.org/ibus-typing-booster/>





# Demo

---



# Thank you!!



- [Anish.developer@gmail.com](mailto:Anish.developer@gmail.com)
- [apatil@redhat.com](mailto:apatil@redhat.com)
- Irc: anish, \_\_anish\_\_

