

# ENTROPY IN NLP

Presented by

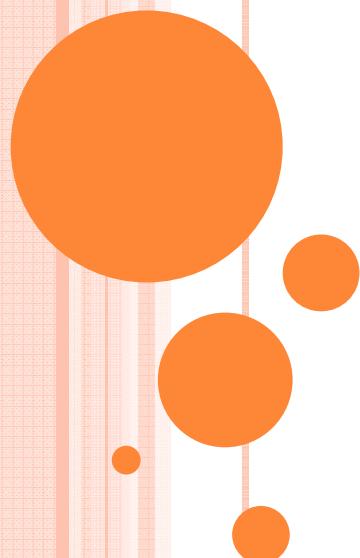
**Avishek Dan (113050011)**

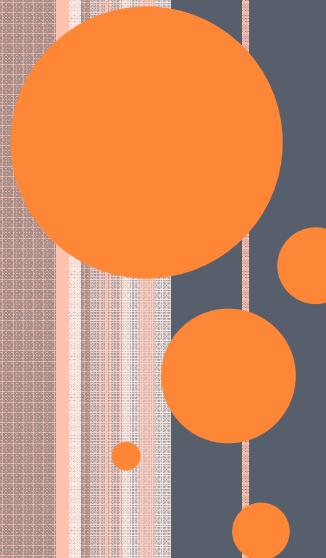
**Lahari Poddar (113050029)**

**Aishwarya Ganesan (113050042)**

Guided by

**Dr. Pushpak Bhattacharyya**





# MOTIVATION



- According to research at an English university, it doesn't matter in what order the letters in a word are, the only important thing is that the first and last letter is at the right place. The rest can be a total mess and you can still read it without a problem. This is because we do not read every letter by itself but the word as a whole.

# 40% REMOVED

Dark s r c res owe neit er side the froze w erwa . T  
trees h d b n tr p d b c t w n o heir white c v ring of  
f s and h y s med e n towa ds e c o her la k and  
ios, i ef l gh . A as ilence reig ed o he  
and. The lan ef a ade l tion, life e s, i  
mo ment, so n nd cold h the pi to it s n t ev n hat  
o sa nes . h e wa hin n it of laughter, bu a lau h r  
et ribe a dn - laught r a was mi hle s as  
the i e of t es a la hte ol as he ro t a d ar ak n  
f th g mne of inf llib i . It wa he mas erful n  
incom un b ew s om of t rnity l ugh n at e lity f fe  
nd e effort f ife. I as t e , savage, fro n-  
h a r hland Wi d.

## 30% REMOVED

Darre forest on either side the frozen waterway. The trees have seen stillness before wide fields white covering of frost, and they stand close towards each other, laden with snow, not even in sight. A vast silence reigns over the land. The landscape was desolation, if one, without men or animals could see that the point of it was not even the sadness. There was a hint in the laughter, up to a laugh overblet many sadnesses a laughe at as worthless as the millions of phosphorus, a ghter cold as the frost and paraking of the grimaces of malability was themselves ruler and community winter laugh fitfully often than the effort of life. It was the Wild, how savage, froze - hearted Northland Wild.

## 20% REMOVED

Dark spruce forest frowned on either side the road. The trees had been stripped by recent wind, their white covering of frost and they seemed to lean towards each other, black and minous, in the fading light. A vast silence reigned over the land. The land itself was desolation, lifeless, without movement, long and cold, having the spirit of silent death upon it. There was a hint of a laugh, but far laughter more terrible than any sense - a laughter that was melancholy as the mile of the sphinx, a laughter colder than the frost and artless than the grimness of infallibility. It was the last and incalculable wisdom of eternity, looking at the futility of life and the fear of life. It was the Wild, the savage, frozen-headed Northland Wild.

# 10% REMOVED

Dark s ru e forest frowned on either side the frozen waterw y. The trees had bee stripped by a recent ind of t eir w ite covering of rost, and they seemed o lean towards each ot er, black and ominous, in the fading li h . A vast silence reigned ver the land The land itself w s a deso ation, lifel ss, without ovement, s l n and cold hat the spiri of it wa not even that of sa ness. here was a hint in it of laughte , but of a lug ter more terrible than any sadness - a ughte that as mirthless s he smile of the sphinx a laug ter cold as the rost a d p rt k ng of the grimness of infal i lity. It as the masterful and incomunica le wisdom of eternity laughing at th futility of life an the effort f life. It was e Wild, the sa ag , froz n- earte Northland Wild.

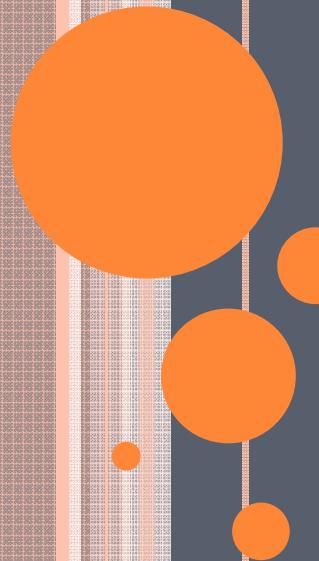
# 0% REMOVED

Dark spruce forest frowned on either side the frozen waterway. The trees had been stripped by a recent wind of their white covering of frost, and they seemed to lean towards each other, black and ominous, in the fading light. A vast silence reigned over the land. The land itself was a desolation, lifeless, without movement, so lone and cold that the spirit of it was not even that of sadness. There was a hint in it of laughter, but of a laughter more terrible than any sadness - a laughter that was mirthless as the smile of the sphinx, a laughter cold as the frost and partaking of the grimness of infallibility. It was the masterful and incommunicable wisdom of eternity laughing at the futility of life and the effort of life. It was the Wild, the savage, frozen-hearted Northland Wild.

From Jack London's "White Fang"

**I**N A previous paper<sup>1</sup> the entropy and redundancy of a language have been defined. The entropy is a statistical parameter which measures in a certain sense, how much information is produced on the average for each letter of a text in the language. If the language is translated into binary digits (0 or 1) in the most efficient way, the entropy  $H$  is the average number of binary digits required per letter of the original language. The redundancy, on the other hand, measures the amount of constraint imposed on a text by the language due to its statistical structure, e.g., in English the high frequency of the letter  $E$ , the strong tendency of  $H$  to follow  $T$  or of  $U$  to follow  $Q$ . It was estimated that when statistical effects extending over not more than eight letters are considered the entropy is roughly 2.3 bits per letter and the redundancy about 50 per cent.

Since then a new method has been found for estimating these quantities which is more sensitive and takes account of long range statistics, influences extending over phrases, sentences, etc. This method is based on a study of the predictability of English; how well can the next letter of a text be predicted when the preceding  $N$  letters are known. The results of some experiments in prediction will be given, and a theoretical analysis of some of the properties of ideal prediction. By combining the experimental and theoretical results it is possible to estimate upper and lower bounds for the entropy and redundancy. From this analysis it appears that in ordinary literar



ENTROPY

# ENTROPY

- Entropy or self-information is the average uncertainty of a single random variable:

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log_2 p(x)$$

- (i)  $H(x) \geq 0$ ,
- (ii)  $H(X) = 0$  only when the value of  $X$  is determinate, hence providing no new information
- From a language perspective, it is the information that is produced on the average for each letter of text in the language

## EXAMPLE: SIMPLE POLYNESIAN

- Random sequence of letters with probabilities:

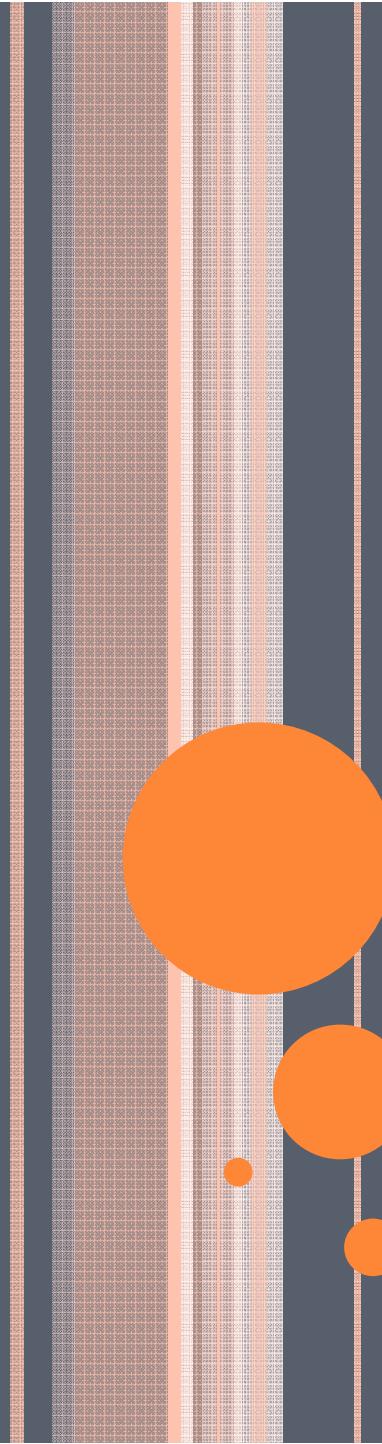
p	t	k	a	i	u
1/8	1/4	1/8	1/4	1/8	1/8

- Per-letter entropy is:

$$\begin{aligned} H(P) &= - \sum_{i \in \{p,t,k,a,i,u\}} P(i) \log P(i) \\ &= -[4 \times \frac{1}{8} \log \frac{1}{8} + 2 \times \frac{1}{4} \log \frac{1}{4}] \\ &= 2\frac{1}{2} \text{ bits} \end{aligned}$$

- Code to represent language:

p	t	k	a	i	u
100	00	101	01	110	111



## SHANON'S EXPERIMENT TO DETERMINE ENTROPY

# SHANNON'S ENTROPY EXPERIMENT

MY PSYCHIC WOULD BE VERY INTERESTED TO LEARN HOW YOU MANAGED TO SWALLOW THAT EGG WHOLE WITHOUT BREAKING IT

7 2 1 2 13 2 1 1 2 1 1 5 3 18 1 1 1 15 2 3 18 2 1 1 1 18 1 1 1 1 1  
1 1 1 1 1 5 1 1 6 3 3 1 1 1 9 3 1 1 5 1 1 1 10 3 1 1 2 1 1 1 1  
1 1 18 15 3 1  
1 1 1 1 1 1 1 1 1 5 1 1 1

The entropy for this experiment is 1.777348

Letters

New Quote

Audio:

On

Off

Source: <http://www.math.ucsd.edu/~crypto/java/ENTROPY/>

# CALCULATION OF ENTROPY

- User as a language model
- Encode number of guesses required
- Apply entropy encoding algorithm (lossless compression)



# ENTROPY OF A LANGUAGE

- Series of approximations

$$F_0, F_1, F_2 \dots F_n$$

$$\begin{aligned} F_N &= - \sum_{i,j} p(b_i, j) \log_2 p_{b_i}(j) \\ &= - \sum_{i,j} p(b_i, j) \log_2 p(b_i, j) + \sum_i p(b_i) \log_2 p(b_i) \end{aligned}$$

$$H = \lim_{n \rightarrow \infty} F_N$$

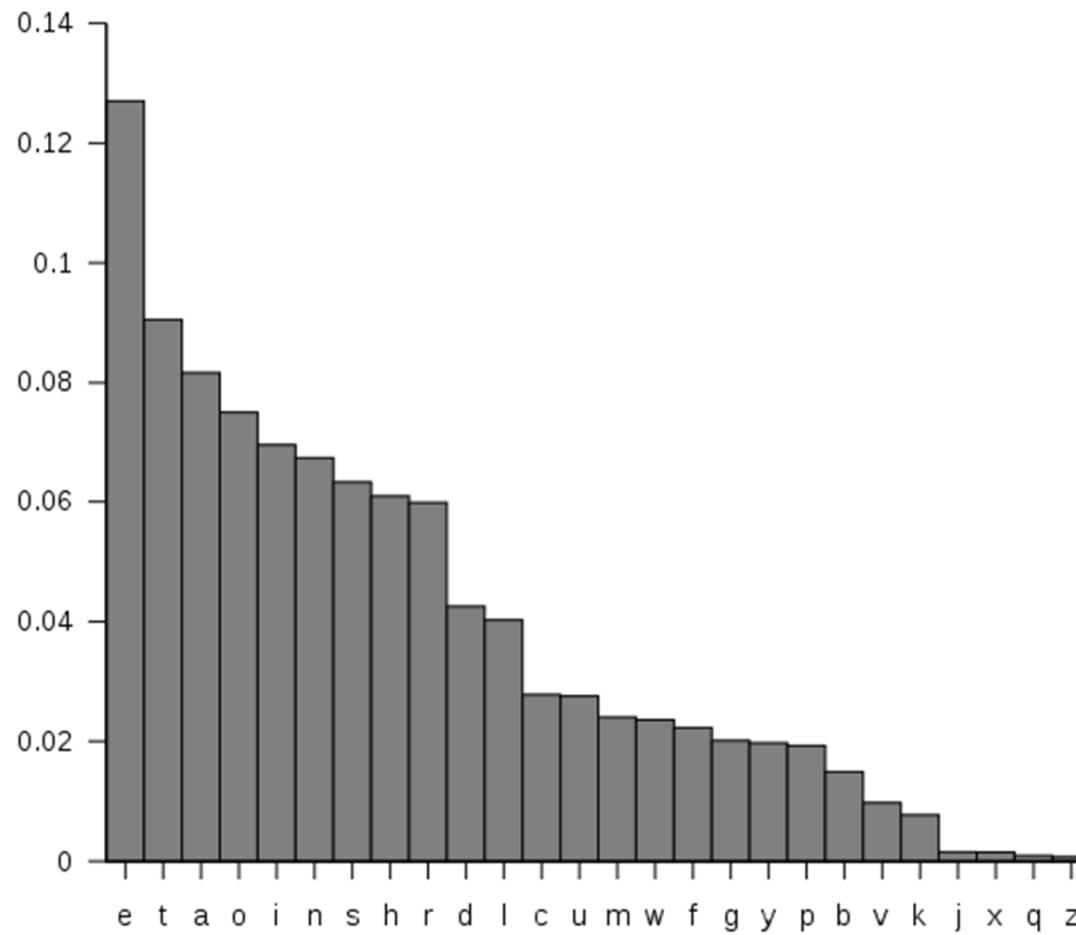


# ENTROPY OF ENGLISH

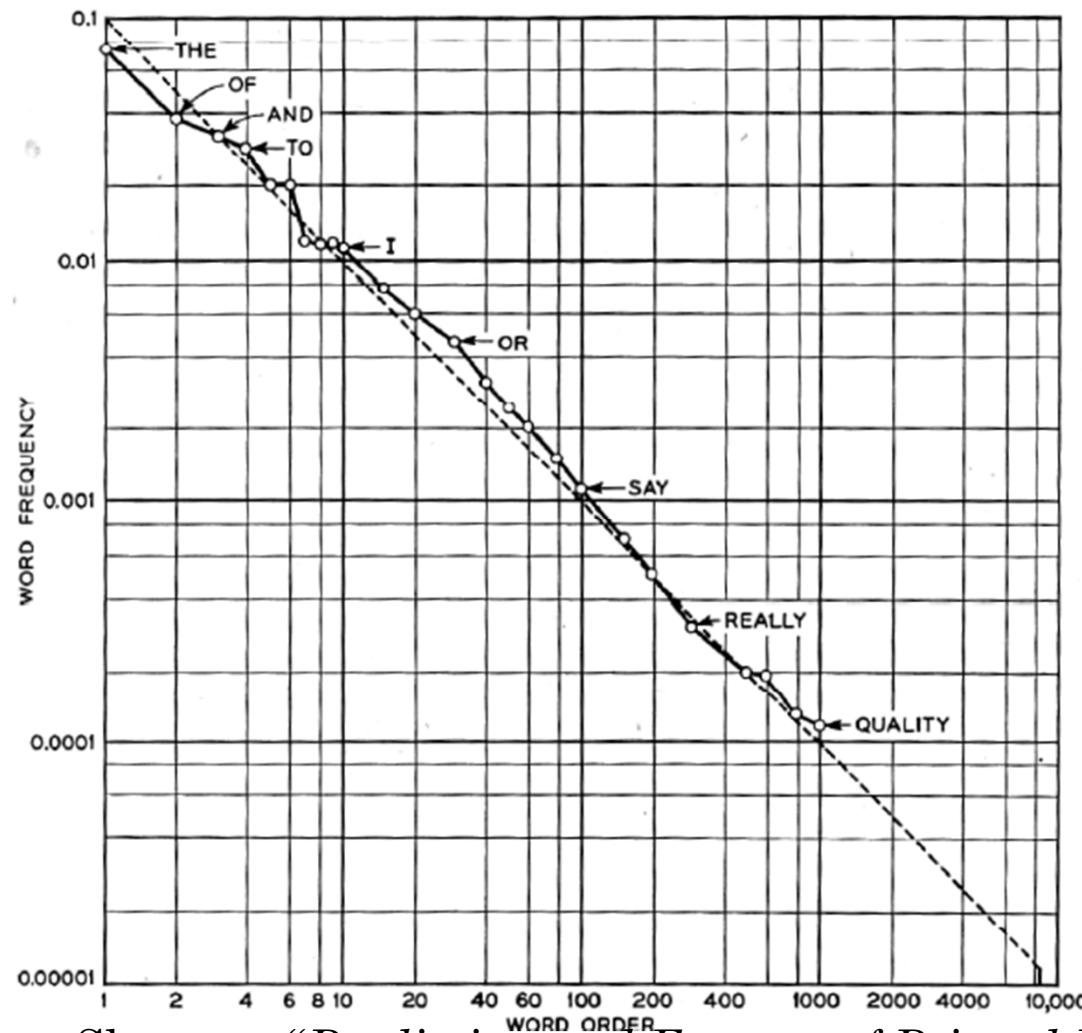
- $F_0 = \log_2 26 = 4.7$  bits per letter
- $F_1 = -\sum_{i=1}^{26} p(i) \log_2 p(i) = 4.14$  bits per letter
- $$\begin{aligned} F_2 &= -\sum_{i,j} p(i, j) \log_2 p_i(j) \\ &= -\sum_{i,j} p(i, j) \log_2 p(i, j) + \sum_i p(i) \log_2 p(i) \\ &= 7.70 - 4.14 = 3.56 \text{ bits per letter.} \end{aligned}$$
- $$\begin{aligned} F_3 &= -\sum_{i,j,k} p(i, j, k) \log_2 p_{ij}(k) \\ &= -\sum_{i,j,k} p(i, j, k) \log_2 p(i, j, k) + \sum_{i,j} p(i, j) \log_2 p(i, j) \\ &\doteq 11.0 - 7.7 = 3.3 \end{aligned}$$



# RELATIVE FREQUENCY OF OCCURRENCE OF ENGLISH ALPHABETS



# WORD ENTROPY OF ENGLISH



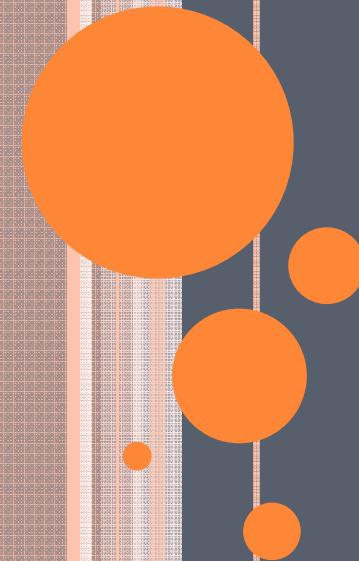
Source: Shannon "Prediction and Entropy of Printed English"

## ZIPF'S LAW

- Zipf's law states that given some corpus of natural language utterances, the frequency of any word is inversely proportional to its rank in the frequency table

$$p_n = \frac{1}{n}$$





# LANGUAGE MODELING

# LANGUAGE MODELING

A language model computes either:

- probability of a sequence of words:

$$P(W) = P(w_1, w_2, w_3, w_4, w_5 \dots w_n)$$

- probability of an upcoming word:

$$P(w_n | w_1, w_2 \dots w_{n-1})$$



# APPLICATIONS OF LANGUAGE MODELS

- POS Tagging
- Machine Translation
  - $P(\text{heavy rains tonight}) > P(\text{weighty rains tonight})$
- Spell Correction
  - $P(\text{about fifteen minutes from}) > P(\text{about fifteen minuets from})$
- Speech Recognition
  - $P(\text{I saw a van}) \gg P(\text{eyes awe of an})$



## N-GRAM MODEL

Chain rule

$$P(w_1 w_2 \dots w_n) = \prod_i P(w_i | w_1 w_2 \dots w_{i-1})$$

Markov Assumption

$$P(w_i | w_1 w_2 \dots w_{i-1}) \approx P(w_i | w_{i-k} \dots w_{i-1})$$

Maximum Likelihood Estimate (for k=1)

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$



# EVALUATION OF LANGUAGE MODELS

A good language model gives a high probability to real English

- Extrinsic Evaluation
  - For comparing models A and B
  - Run applications like POSTagging, translation in each model and get an accuracy for A and for B
  - Compare accuracy for A and B
- Intrinsic Evaluation
  - Use of cross entropy and perplexity
  - True model for data has the lowest possible entropy / perplexity



## RELATIVE ENTROPY

- For two probability mass functions,  $p(x)$ ,  $q(x)$  their relative entropy:

$$D(p \parallel q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)}$$

- Also known as *KL divergence*, a measure of how different two distributions are.

$$D(p \parallel q) = E_p \left( \log \frac{p(X)}{q(X)} \right)$$



## CROSS ENTROPY

- Entropy as a measure of how surprised we are, measured by pointwise entropy for model m:

$$H(w/h) = -\log(m(w/h))$$

- Produce q of real distribution to minimize  $D(p \parallel q)$
- The cross entropy between a random variable X with  $p(x)$  and another q (a model of p) is:

$$\begin{aligned} H(X, q) &= H(X) + D(p \parallel q) \\ &= - \sum_x p(x) \log q(x) \end{aligned}$$

$$H(L, m) \approx -\frac{1}{n} \log m(x_{1:n})$$



## PERPLEXITY

- Perplexity is defined as

$$PP(x_{1:n}, m) = 2^{H(x_{1:n}, m)}$$

- Probability of the test set assigned by the language model, normalized by the number of word

$$PP(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_{i-1})}}$$

- Most common measure to evaluate language models



# EXAMPLE

Training set

*there is a big house  
i buy a house  
they buy the new house*

Model

$$\begin{array}{lll} p(\text{big}|a) = 0.5 & p(\text{is}|\text{there}) = 1 & p(\text{buy}|\text{they}) = 1 \\ p(\text{house}|a) = 0.5 & p(\text{buy}|i) = 1 & p(a|\text{buy}) = 0.5 \\ p(\text{new}|\text{the}) = 1 & p(\text{house}|\text{big}) = 1 & p(\text{the}|\text{buy}) = 0.5 \\ p(a|\text{is}) = 1 & p(\text{house}|\text{new}) = 1 & p(\text{they}|<\text{s}>) = .333 \end{array}$$

Test sentence  $S$ : *they buy a big house*

$$p(S) = \underbrace{0.333}_{\text{they}} \times \underbrace{1}_{\text{buy}} \times \underbrace{0.5}_{\text{a}} \times \underbrace{0.5}_{\text{big}} \times \underbrace{1}_{\text{house}} = 0.0833$$



## EXAMPLE (CONTD.)

- Cross Entropy:

$$\begin{aligned} H(p, m) &= -\frac{1}{5} \log p(S) \\ &= -\frac{1}{5} (\underbrace{\log 0.333}_{they} + \underbrace{\log 1}_{buy} + \underbrace{\log 0.5}_{a} + \underbrace{\log 0.5}_{big} + \underbrace{\log 1}_{house}) \\ &= -\frac{1}{5} (-1.586 + 0 + -1 + -1 + 0) = 0.7173 \end{aligned}$$

- Perplexity:

$$PP = 1.6441$$



## SMOOTHING

- Bigrams with zero probability - cannot compute perplexity
- When we have sparse statistics Steal probability mass to generalize better ones.
- Many techniques available
- Add-one estimation : add one to all the counts

$$P_{Add-1}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i) + 1}{c(w_{i-1}) + V}$$

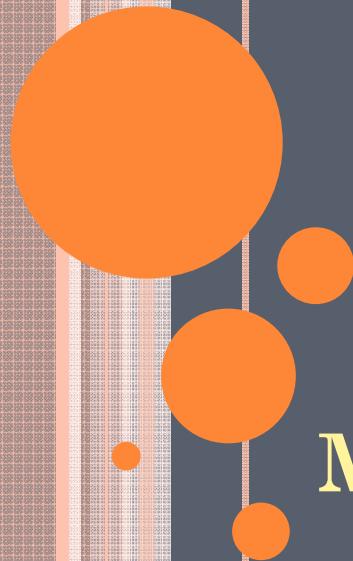


## PERPLEXITY FOR N-GRAM MODELS

- Perplexity values yielded by n-gram models on English text range from about 50 to almost 1000 (corresponding to cross entropies from about 6 to 10 bits/word)
- Training: 38 million words from WSJ by Jurafsky
- Vocabulary: 19,979 words
- Test: 1.5 million words from WSJ

N-gram Order	Unigram	Bigram	Trigram
Perplexity	962	170	109





MAXIMUM ENTROPY MODEL

# STATISTICAL MODELING

- Constructs a stochastic model to predict the behavior of a random process
  - Given a sample, represent a process
  - Use the representation to make predictions about the future behavior of the process
- 
- Eg: Team selectors employ batting averages, compiled from history of players, to gauge the likelihood that a player will succeed in the next match. Thus informed, they manipulate their lineups accordingly



# STAGES OF STATISTICAL MODELING

1. Feature Selection :Determine a set of statistics that captures the behavior of a random process.
  
2. Model Selection: Design an accurate model of the process--a model capable of predicting the future output of the process.



## MOTIVATING EXAMPLE

- Model an expert translator's decisions concerning the proper French rendering of the English word *on*
- A model( $p$ ) of the expert's decisions assigns to each French word or phrase( $f$ ) an estimate,  $p(f)$ , of the probability that the expert would choose  $f$  as a translation of *on*
- Our goal is to
  - Extract a set of facts about the decision-making process from the sample
  - Construct a model of this process



## MOTIVATING EXAMPLE

- A clue from the sample is **the list of allowed translations**
  - $on \rightarrow \{sur, dans, par, au bord de\}$
- With this information in hand, we can impose our first constraint on  $p$ :

$$p(sur) + p(dans) + p(par) + p(au bord de) = 1$$

- The most uniform model will divide the probability values equally
- Suppose we notice that the expert chose either *dans* or *sur* 30% of the time, then a second constraint can be added

$$p(dans) + p(sur) = 3/10$$

- Intuitive Principle: Model all that is known and assume nothing about that which is unknown



# MAXIMUM ENTROPY MODEL

- A random process which produces an output value  $y$ , a member of a finite set  $\mathbf{Y}$ .
  - $y \in \{\text{sur, dans, par, au bord de}\}$
- The process may be influenced by some contextual information  $x$ , a member of a finite set  $\mathbf{X}$ .
  - $x$  could include the words in the English sentence surrounding *on*
- A stochastic model: accurately representing the behavior of the random process
  - Given a context  $x$ , the process will output  $y$



# MAXIMUM ENTROPY MODEL

- Empirical Probability Distribution:

$$\tilde{p}(x, y) \equiv \frac{1}{N} \times \text{number of times that } (x, y) \text{ occurs in the sample}$$

- Feature Function:

$$f(x, y) = \begin{cases} 1 & \text{if } y = \text{sur and hold follows on} \\ 0 & \text{otherwise} \end{cases}$$

- Expected Value of the Feature Function

- For training data:

$$\tilde{p}(f) \equiv \sum_{x, y} \tilde{p}(x, y) f(x, y) \quad (1)$$

- For model:

$$p(f) \equiv \sum_{x, y} \tilde{p}(x) p(y | x) f(x, y) \quad (2)$$



## MAXIMUM ENTROPY MODEL

- To accord the model with the statistic, the expected values are constrained to be equal

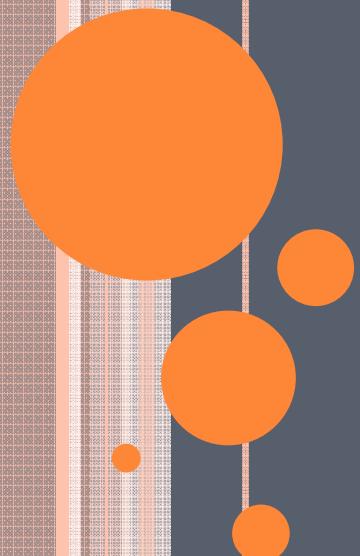
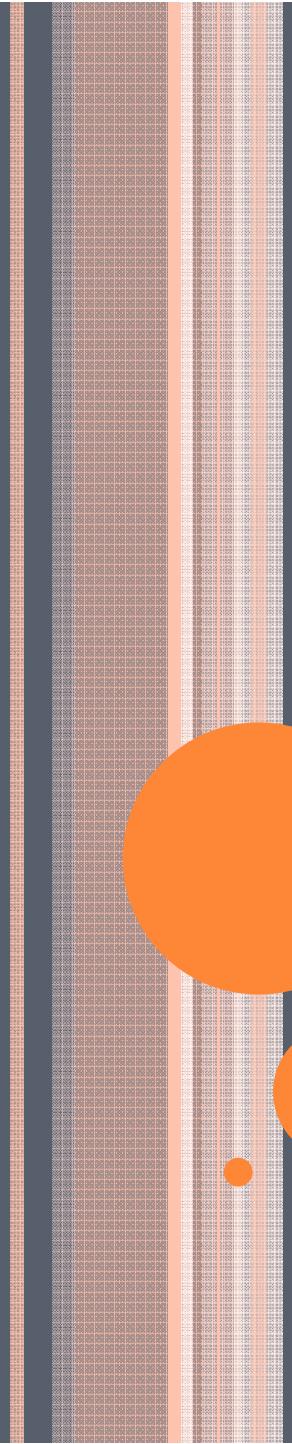
$$p(f) = \tilde{p}(f) \quad (3)$$

- Given  $n$  feature functions  $f_i$ , the model  $p$  should lie in the subset  $\mathcal{C}$  of  $\mathbf{P}$  defined by
- Choose the model  $p^*$  with maximum entropy  $H(p)$ :

$$\mathcal{C} \equiv \{p \in \mathbf{P} \mid p(f_i) = \tilde{p}(f_i) \text{ for } i \in \{1, 2, \dots, n\}\} \quad (4)$$

$$H(p) = - \sum_{x,y} \tilde{p}(x) p(y|x) \log p(y|x) \quad (5)$$





# APPLICATIONS OF MEM

# STATISTICAL MACHINE TRANSLATION

- A French sentence  $F$ , is translated to an English sentence  $E$  as:

$$E = \arg \max_E p(F | E) P(E)$$

- Addition of MEM can introduce context-dependency:
  - $P_e(f | x)$  : Probability of choosing  $e$  (English) as the rendering of  $f$  (French) given the context  $x$

Translation	$e_{-3}$	$e_{-2}$	$e_{-1}$	$e_{+1}$	$e_{+2}$	$e_{+3}$
<i>dans</i>	<i>the</i>	<i>committee</i>	<i>stated</i>	.	<i>a</i>	<i>letter</i>
<i>à</i>	<i>work</i>	<i>was</i>	<i>required</i>	.	<i>respect</i>	<i>to</i>
<i>au cours de</i>				.	<i>the</i>	<i>the</i>
<i>dans</i>	<i>by</i>	<i>the</i>	<i>government</i>	.	<i>fiscal</i>	<i>year</i>
<i>à</i>	<i>of</i>	<i>diphtheria</i>	<i>reported</i>	.	<i>the</i>	<i>same</i>
<i>de</i>	<i>not</i>	<i>given</i>	<i>notice</i>	.	<i>Canada</i>	<i>postal</i>
				.	<i>, the</i>	<i>by</i>
				.	<i>ordinary</i>	<i>way</i>

# PART OF SPEECH TAGGING

- The probability model is defined over  $H \times T$ 
  - $H$  : set of possible word and tag contexts(histories)
  - $T$  : set of allowable tags
- Entropy of the distribution :

$$H(p) = - \sum_{h \in H, t \in T} p(h, t) \log p(h, t)$$

- Sample Feature Set :

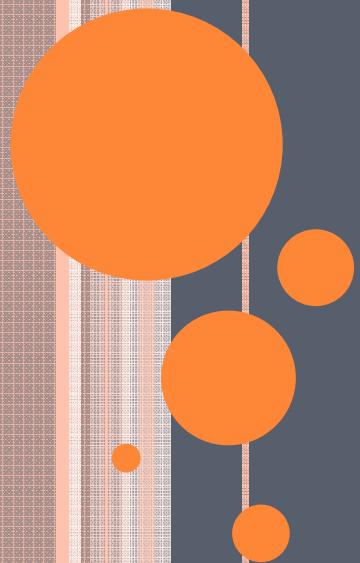
Condition	Features	
$w_i$ is not rare	$w_i = X$	$\& t_i = T$
$w_i$ is rare	$X$ is prefix of $w_i$ , $ X  \leq 4$	$\& t_i = T$
	$X$ is suffix of $w_i$ , $ X  \leq 4$	$\& t_i = T$
	$w_i$ contains number	$\& t_i = T$
	$w_i$ contains uppercase character	$\& t_i = T$
	$w_i$ contains hyphen	$\& t_i = T$
$\forall w_i$	$t_{i-1} = X$	$\& t_i = T$
	$t_{i-2}t_{i-1} = XY$	$\& t_i = T$
	$w_{i-1} = X$	$\& t_i = T$
	$w_{i-2} = X$	$\& t_i = T$
	$w_{i+1} = X$	$\& t_i = T$
	$w_{i+2} = X$	$\& t_i = T$

- Precision : 96.6%

# PREPOSITION PHRASE ATTACHMENT

- MEM produces a probability distribution for the PP-attachment decision using only information from the verb phrase in which the attachment occurs
- conditional probability of an attachment is  $p(d | h)$ 
  - $h$  is the history
  - $d \in \{0, 1\}$ , corresponds to a noun or verb attachment (respectively)
- Features: Testing for features should only involve
  - Head Verb (V)
  - Head Noun (N1)
  - Head Preposition (P)
  - Head Noun of the Object of the Preposition (N2)
- Performance :
  - Decision Tree : 79.5 %
  - MEM : 82.2%





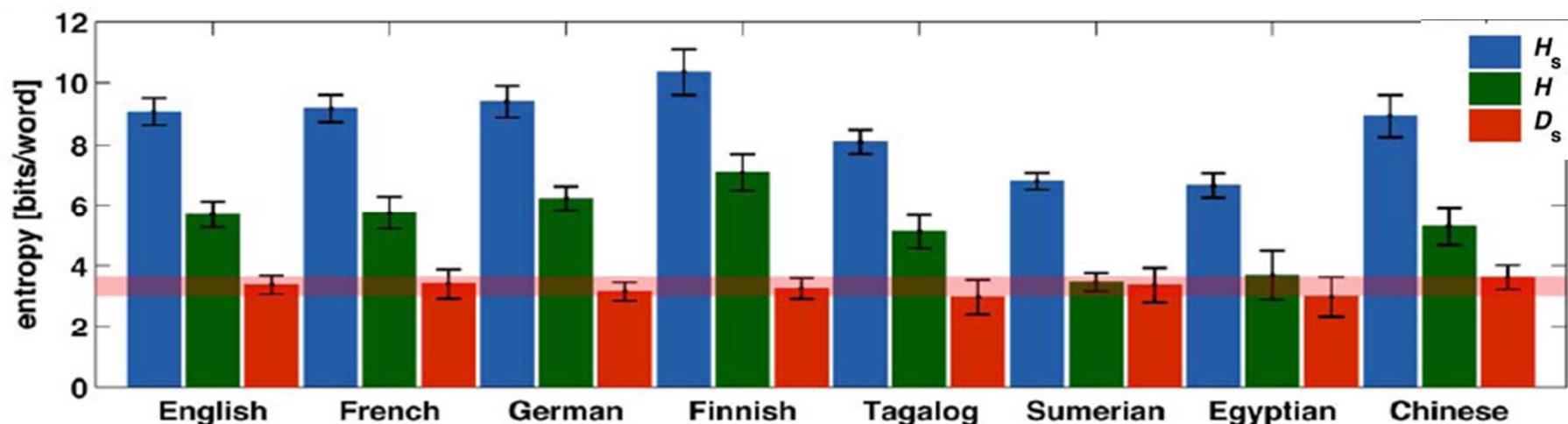
# ENTROPY OF OTHER LANGUAGES

# ENTROPY OF EIGHT LANGUAGES BELONGING TO FIVE LINGUISTIC FAMILIES

- Indo-European: English, French, and German
- Finno-Ugric: Finnish
- Austronesian: Tagalog
- Isolate: Sumerian
- Afroasiatic: Old Egyptian
- Sino-Tibetan: Chinese

$H_s$ - entropy when words are random ,  $H$ - entropy when words are ordered

$$D_s = H_s - H$$



Source: Universal Entropy of Word Ordering Across Linguistic Families

# ENTROPY OF HINDI

- *Zero-order* : 5.61 bits/symbol.
- *First-order* : 4.901 bits/symbol.
- *Second-order* : 3.79 bits/symbol.
- *Third-order* : 2.89 bits/symbol.
- *Fourth-order* : 2.31 bits/symbol.
- *Fifth-order* : 1.89 bits/symbol.

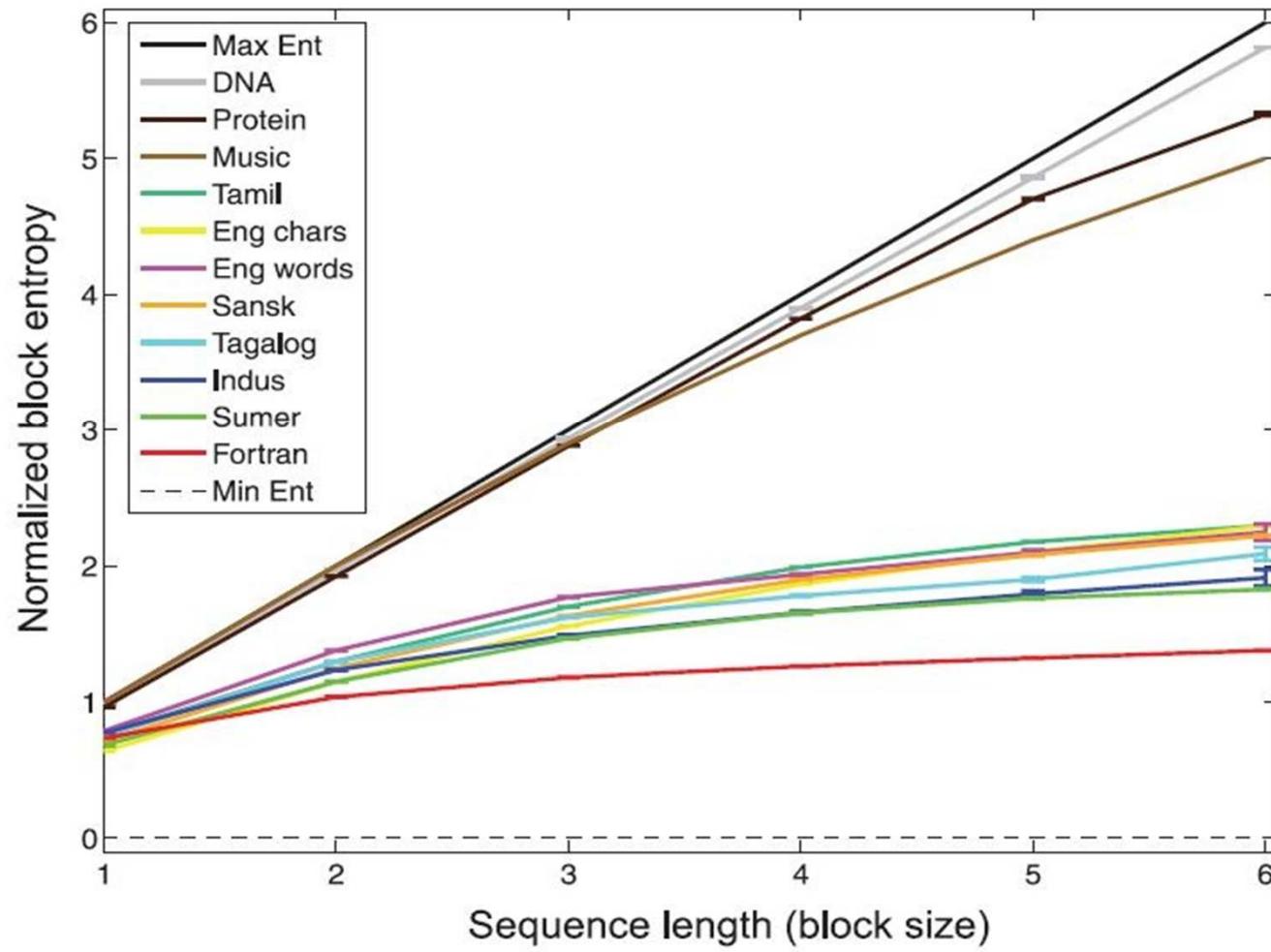


# ENTROPY TO PROVE THAT A SCRIPT REPRESENT LANGUAGE

- Pictish (a Scottish, Iron Age culture) symbols revealed as a written language through application of Shannon entropy
- In Entropy, the Indus Script, and Language by proving the block entropies of the Indus texts remain close to those of a variety of natural languages and far from the entropies for unordered and rigidly ordered sequences



# ENTROPY OF LINGUISTIC AND NON-LINGUISTIC LANGUAGES



Source: *Entropy, the Indus Script, and Language*

# CONCLUSION

- The concept of entropy dates back to biblical times but it has wide range of applications in modern NLP
- NLP is heavily supported by Entropy models in various stages. We have tried to touch upon few aspects of it.



## REFERENCES

- C. E. Shannon, *Prediction and Entropy of Printed English*, Bell System Technical Journal Search, Volume 30, Issue 1, January 1951
- <http://www.math.ucsd.edu/~crypto/java/ENTROPY/>
- Chris Manning and Hinrich Schütze, Foundations of Statistical Natural Language Processing, MIT Press. Cambridge, MA: May 1999.
- Daniel Jurafsky and James H. Martin, SPEECH and LANGUAGE PROCESSING An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition Second Edition ,January 6, 2009
- Marcelo A. Montemurro, Damian H. Zanette, *Universal Entropy of Word Ordering Across Linguistic Families*, PLoS ONE 6(5): e19875. doi:10.1371/journal.pone.0019875

- Rajesh P. N. Rao, Nisha Yadav, Mayank N. Vahia, Hrishikesh Joglekar, Ronojoy Adhikari, Iravatham Mahadevan, *Entropy, the Indus Script, and Language: A Reply to R. Sproat*
- Adwait Ratnaparkhi, *A Maximum Entropy Model for Prepositional Phrase Attachment*, HLT '94
- Berger, A Maximum Entropy Approach to Natural Language Processing, 1996
- Adwait Ratnaparkhi, A Maximum Entropy Model for Part-Of-Speech Tagging, 1996

