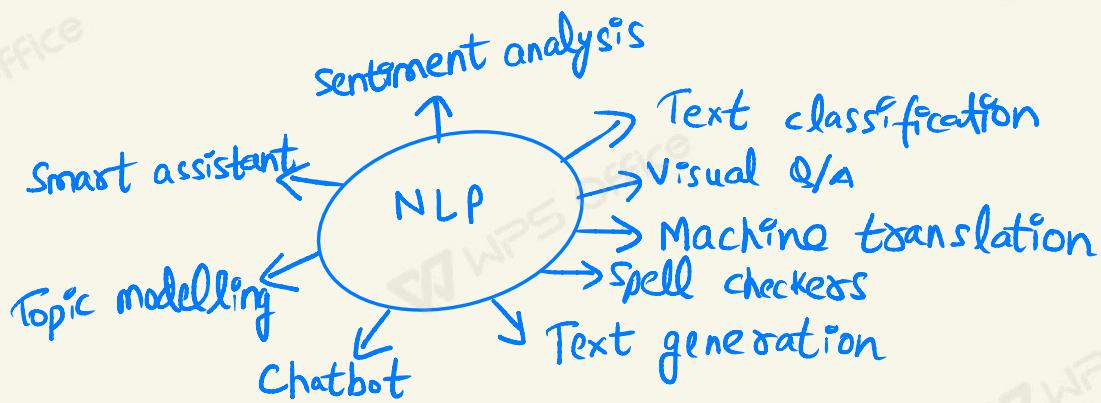


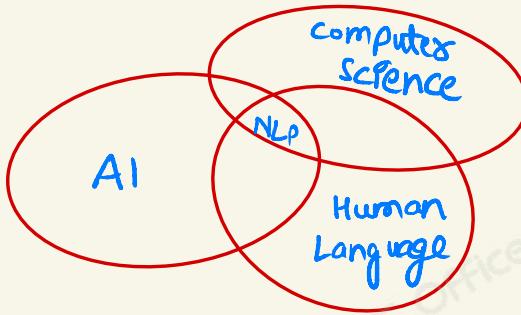


# Application of NLP



## Best NLP language models

- ① BERT → Bidirectional Encoder representation from Transformation
- ② ROBERTa → Robustly optimized BERT pre-trained approach
- ③ XLNET
- ④ Open AI's GPT 3
- ⑤ ALBERT
- ⑥ T5



→ NLP is a branch of AI which deals with communication between human language and computer that allows machine/computer to understand human language.

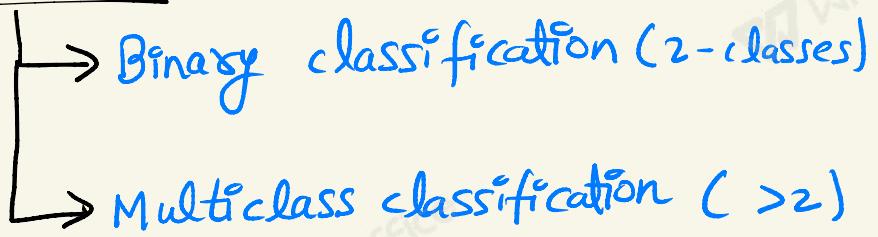
## Most popular Libraries of NLP

- ① Hugging face (HF) Transformers
- ② Spacy
- ③ NLTK
- ④ Gensim
- ⑤ Fairseq

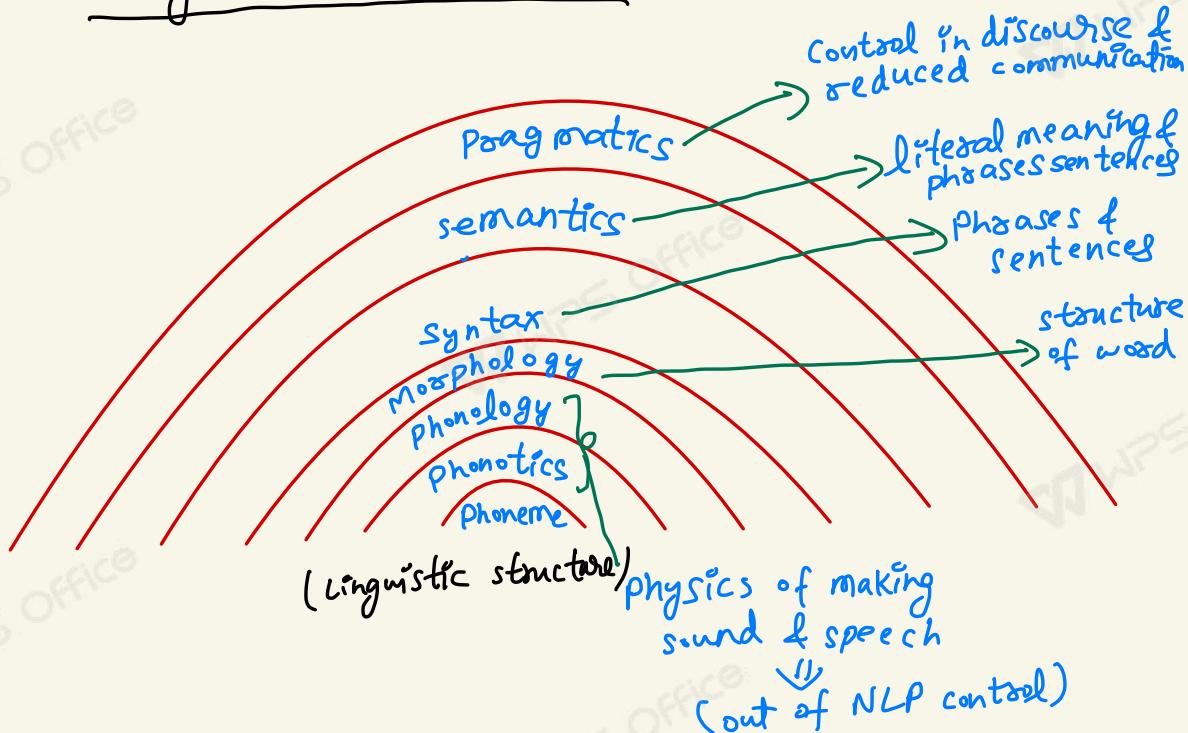
~~HW~~      zero-shot  
                few-shot

{ self supervised  
 { unsupervised  
 { semi-supervised

## Classification

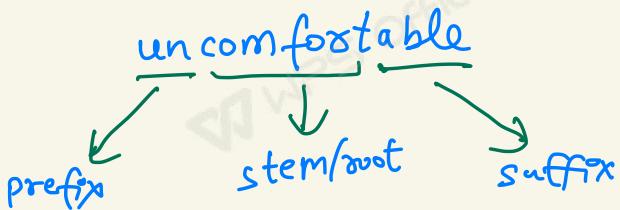
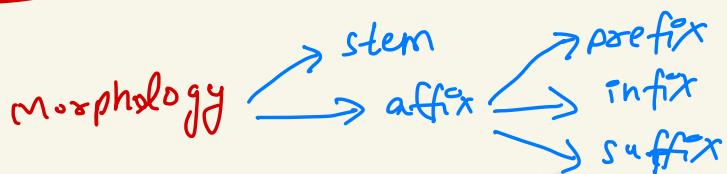


# Linguistic structure



morphology → study of structure of word.

Lexicography → study of identifying words & phrases



## (1) Inflection

## (2) Derivation

## (3) Compounding

stem → plural, past, progressive

Ex:- jump → jumps, jumped, jumping  
like → linked, linked, liking

## (2) Derivation

paint + er → painter

re + paint → reprint

paint + able → printable

## (3) Compounding

corpora → corpus

  
text data

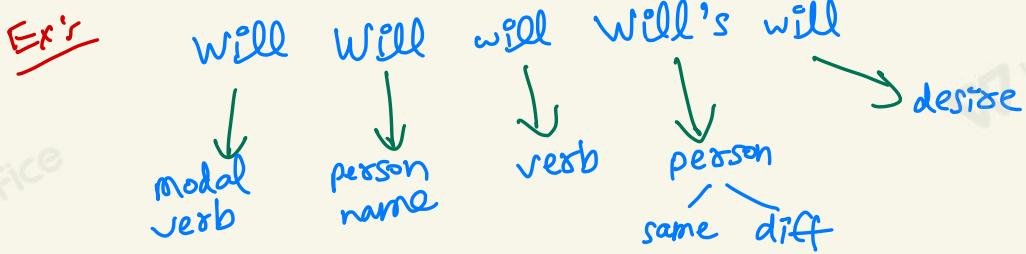
### stem + stem

house + boat → houseboat

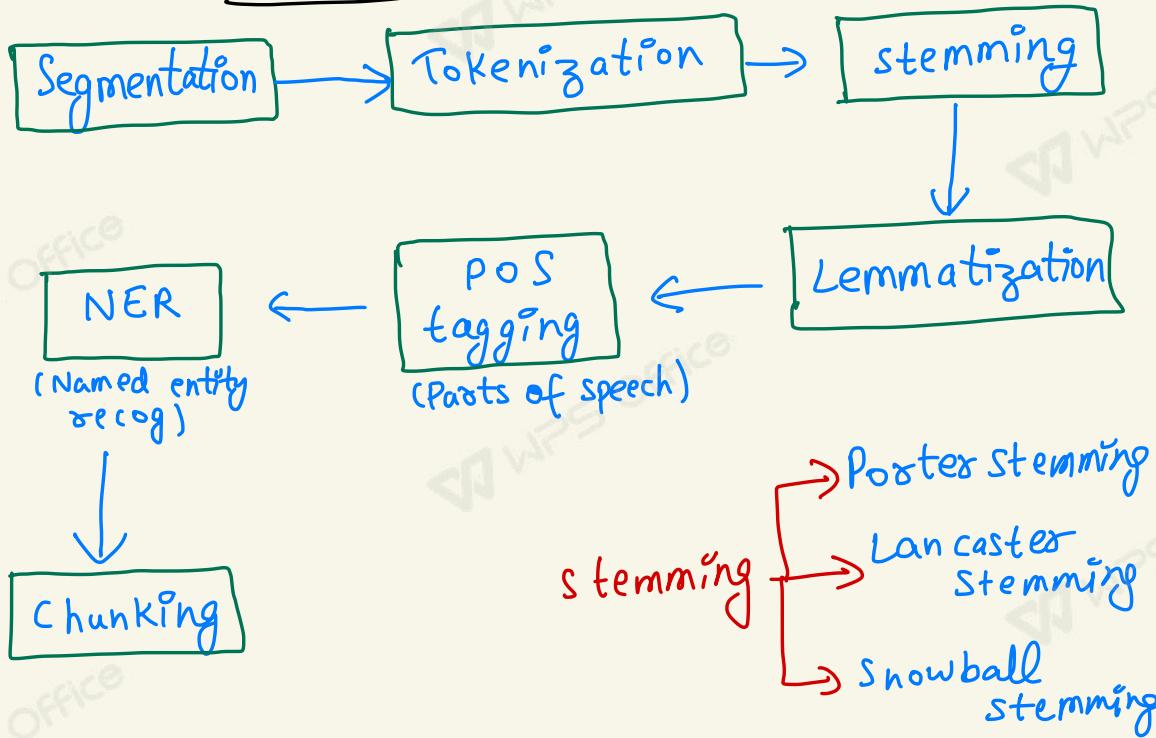
boat + house → boathouse

(order matters)

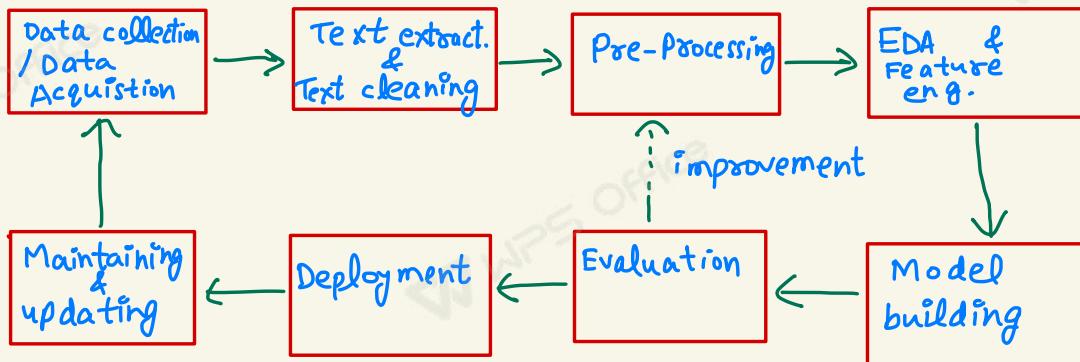
morphological rich → order does not matter  
(same meaning)



## NLP stages / NLP Pipeline



# NLP PIPELINE (In widen properties)/Lifecycle



# Probability Language Model

- ↳ Joint probability
- ↳ conditional probability
- ↳ marginal probability

$$P(A, B) = P(A) \cdot P(B|A)$$

$$P(A, B, C) = P(A) P(B|A) P(C|A, B)$$

$$P(A, B, C, D) = P(A) \cdot P(B|A) \cdot P(C|A, B) \cdot P(D|A, B, C)$$

## Markov's Assumption

- n-gram
- uni-gram
- bi-gram
- tri-gram  $P(w_i | w_{i-1})$
- $P(w_i | w_{i-2} w_{i-1})$

## Probability Overview: Basics of Probability Theory

tossing two coins simultaneously,

possible outcomes  $\Rightarrow$

$$S = \{TT, TH, HT, HH\}$$

1. Sample space
2. Random experiment
3. Favourable event
4. Success.
5. Random variable

# Probability Distribution

## Discrete

1. General Discrete
2. Binomial / Bernoulli's
3. Hyper geometric
4. Geometric
5. Poisson's

## Continuous

1. General continuous
2. Uniform distribution
3. Exponential distribution
4. Normal / Gaussian
5. Standard normal

## Random variable

DRV

↓  
PMF

(Prob mass func<sup>n</sup>)

CRV

↓  
pdf

(Prob density func<sup>n</sup>)

RV = X :	0	1	2
$f(n) = p(x)$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$
$F(n)$	$\frac{1}{4}$	$\frac{1}{2} + \frac{1}{4}$	1

$$\sum_n f(n) = 1$$

$$\sum_n p(n) = 1$$

(1) Expectation  $E(n)$

$$\text{mean of } X = \mu = \sum_n n f(n)$$

Avg. of  $X$

(2) Variance

$$\text{Var}(X) = \frac{1}{N} \sum (n - \mu)^2 = E(X^2) - E(X)^2$$

(3) Standard deviation ( $\sigma$ ) =  $\sqrt{\text{Var}(X)}$

Q. Consider the following PMF of a random variable

$n$ ,

$$P(n, q) = \begin{cases} q & \text{if } n=0 \\ \cdot & \cdot \\ 1-q & \text{if } n=1 \\ 0 & \text{otherwise} \end{cases}$$

if  $q=0.4$ , variance?

Q. A machine produces 0, 1 or 2 defective pieces in a day with associated prob of  $\frac{1}{6}, \frac{2}{3}$  and  $\frac{1}{6}$  respectively. Then mean value & the variance of the no. of defective pieces produced by the machine.

(a)  $1, \frac{1}{3}$       (c)  $1, \frac{4}{3}$

(b)  $\frac{1}{3}, 1$       (d)  $\frac{4}{3}, \frac{1}{3}$

Ans 1:

X	0	1
$f(n)$	0.4	0.6

$$\begin{aligned} E(X) &= 0 \times 0.4 + 1 \times 0.6 \\ &= 0.6 \end{aligned}$$

$$\begin{aligned} \text{var}(X) &= E(X^2) - E(X)^2 \\ &= (0 \times 0.4 + 1 \times 0.6) - 0.36 \\ &= 0.6 - 0.36 \\ &= 0.24 \end{aligned}$$

Ans 2:  $E(x) = \sum n f(n)$

$$\text{var}(x) = \left( 1 \times \frac{2}{3} + 2 \times \frac{1}{3} - 1 \right) = \frac{4}{3} - 1 = \frac{1}{3}$$

Q. The function  $P(n)$  is given by  $P(n) = \frac{A}{n^k}$  where  $A$  and  $k$  are constants with  $k > 1$  and  $1 \leq n < \infty$  and  $P(n) = 0$  for  $-\infty < n < 1$  for  $P(n)$  to be a probability density function the value of  $A$  should be equal to

- (A)  $k-1$       (B)  $k+1$       (C)  $1/k-1$       (D)  $1/k+1$

$$\Rightarrow \sum P(n) = 1$$

$$\Rightarrow \sum_{n=1}^{\infty} \frac{A}{n^k} = 1$$

$$\Rightarrow A \left[ \frac{1}{1^k} + \frac{1}{2^k} + \frac{1}{3^k} + \dots \right] = 1$$

$$\Rightarrow \int_{n=1}^{\infty} Ax^k dx = 1$$

$$\Rightarrow Ax \left[ \frac{x^{-k+1}}{-k+1} \right]_{n=1}^{\infty} = 1$$

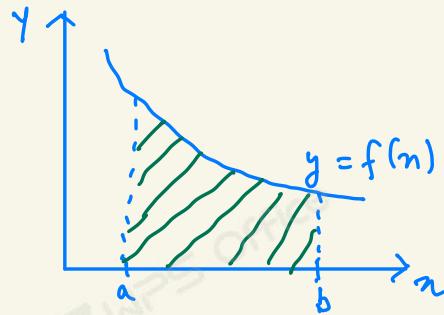
$$\Rightarrow A \left[ \frac{1}{n^{-k} (-k+1)} \right]_{n=1}^{\infty} = 1$$

$$\Rightarrow Ax \left[ 0 - \frac{1}{-k+1} \right] = 1$$

$$\Rightarrow \frac{A}{k-1} = 1$$

$$\Rightarrow A = k-1$$

# General Continuous Probability Distribution



$x$  is a continuous random variable.

$$P(x=a) = 0$$

$$P(x=b) = 0$$

$$P(a < x < b) = \int_a^b f(x) dx$$

(i)  $x \rightarrow c.r.v$  then,  
probability at a single point = 0  
 $P(x=a) = 0$ .

(ii)  $x \rightarrow c.r.v. \Rightarrow P(a \leq x \leq b) = \int_a^b f(x) dx$

(iii)  $P(-\infty < x < \infty) = \text{Area under the density func^n curve.}$   
 $= \int_{-\infty}^{\infty} f(x) dx$

(iv)  $f(x) \geq 0$ , always positive

Q. If  $n$  is a c.r.v. having pdf is given by

$$f(n) = \begin{cases} cn^2, & 0 \leq n < 1 \\ cn, & 1 \leq n < 2 \\ 0, & \text{otherwise} \end{cases}$$

(i) Find  $c$

$$(\text{ii}) \quad P(1/2 < n < 3/2)$$

$$(\text{iii}) \quad P(n < 3/2)$$

$$(\text{iv}) \quad P(n \geq 1/2)$$

$$(i) \quad \int_{-\infty}^{\infty} f(n) dn = 1$$

$$\Rightarrow \int_0^1 cn^2 \cdot dn + \int_1^2 cn \cdot dn = 1$$

$$\Rightarrow \frac{cn^3}{3} \Big|_0^1 + \frac{cn^2}{2} \Big|_1^2 = 1$$

$$\Rightarrow \frac{c}{3} + \left[ c \times 2 - \frac{c}{2} \right] = 1$$

$$\Rightarrow \frac{c}{3} + \frac{3c}{2} = 1$$

$$\Rightarrow \frac{2c + 9c}{6} = 1$$

$$\Rightarrow c = 6/11$$

$$\begin{aligned}
 \text{(ii)} \quad p(n) &= \int_{1/2}^1 cn^2 \cdot dn + \int_1^{3/2} cn \cdot dn \\
 &= \left. \frac{cn^3}{3} \right|_{1/2}^1 + \left. \frac{cn^2}{2} \right|_1^{3/2} \\
 &= \frac{c}{3} - \frac{c}{24} + \frac{cx9}{8} - \frac{c}{2} \\
 &= \frac{6}{11} \left[ \frac{1}{3} - \frac{1}{24} + \frac{9}{8} - \frac{1}{2} \right] \\
 &= \frac{6}{11} \times \frac{8 - 1 + 27 - 12}{24} \\
 &= \frac{6}{11} \times \frac{22}{24} = \frac{1}{2}
 \end{aligned}$$

$$\begin{aligned}
 \text{(iii)} \quad p &= \int_0^1 cn^2 \cdot dn + \int_1^{3/2} cn \cdot dn \\
 &= \left. \frac{cn^3}{3} \right|_0^1 + \left. \frac{cn^2}{2} \right|_1^{3/2} \\
 &= \frac{c}{3} + \frac{cx9}{8} - \frac{c}{2} \\
 &= \frac{6}{11} \times \frac{8 + 27 - 12}{24} \\
 &= \frac{6}{11} \times \frac{23}{24} = \frac{23}{44}
 \end{aligned}$$

$$\begin{aligned}
 \text{(i) } P(X \geq 1/2) &= \int_{1/2}^{\infty} cn^2 dn + \int_1^{\infty} cn dn \\
 &= \left[ \frac{cn^3}{3} \right]_{1/2}^1 + \left[ \frac{cn^2}{2} \right]_1^{\infty} \\
 &= \frac{c}{3} - \frac{c}{24} + \frac{c \cdot 4}{2} - \frac{c}{2} \\
 &= \frac{6}{11} \left[ \frac{1}{3} - \frac{1}{24} + \frac{4}{2} - \frac{1}{2} \right] \\
 &= \frac{6}{11} \times \frac{8 - 1 + 36}{24} \\
 &= \frac{43}{4} \times \frac{1}{11}
 \end{aligned}$$

Q. For the function  $f(n) = a + bn$ ,  $0 < n < 1$  to be a valid P.d.f. which one of the following statement is true.

- |                |                  |
|----------------|------------------|
| (a) $a=1, b=4$ | (b) $a=0.5, b=1$ |
| (c) $a=0, b=1$ | (d) $a=1, b=-1$  |

$$\begin{aligned}
 \Rightarrow \int_0^1 (a+bn) dn &= 1 \\
 \Rightarrow \left[ an + \frac{bn^2}{2} \right]_0^1 &= 1
 \end{aligned}$$

$$a + \frac{b}{2} = 1$$

(b)

Q. Find the value of  $h$ ,



$$\Rightarrow A_1 + A_2 + A_3 = 1$$

$$\Rightarrow \frac{1}{2} \times 1 \times h + \frac{1}{2} \times 1 \times 2h + \frac{1}{2} \times 1 \times 3h = 1$$

$$\Rightarrow h = \frac{1}{3}$$

Q.

$$f(x) = \begin{cases} 0 & , -\infty < x < 0 \\ \frac{2}{11}x^3 & , 0 \leq x < 1 \\ \frac{3x^2-1}{11} & , 1 \leq x < 2 \\ 1 & , 2 \leq x < \infty \end{cases}$$

(i) Find  $P(\frac{1}{2} < x < \frac{3}{2})$

(ii)  $P(x \geq \frac{1}{2})$

(iii)  $P(x < \frac{3}{2})$

Q1 If  $X$  is uniformly distributed in  $(0, 10)$  then find

- (i)  $f(n)$
- (ii) mean, variance, std deviation
- (iii)  $P(2 < X < 6)$
- (iv)  $P(0 < X < 5)$
- (v)  $P(X \leq 3)$
- (vi)  $P(X > 8)$

Q2 If a random variable is uniformly distributed with mean 1 and variance  $\frac{1}{3}$ , then  $P(X < \frac{1}{2}) = ?$

Q3 The p.d.f of a random variable  $X$  is  
$$f(x) = \begin{cases} \frac{x}{4}(4-x^2) & \text{for } 0 \leq x \leq 2 \\ 0 & \text{otherwise} \end{cases}$$

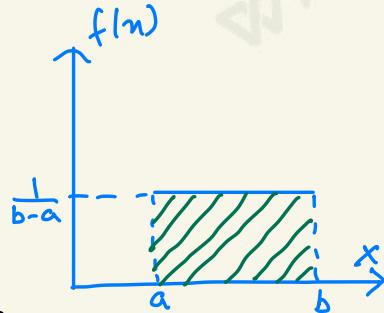
Find the mean.

### Uniform Density

If  $X$  is uniformly distributed continuous random variable, then the probability density function in the finite interval  $[a, b]$  is given as

$$f(n) = \begin{cases} \frac{1}{b-a} & , a \leq x \leq b \\ 0 & , \text{otherwise} \end{cases}$$

$$f(n) = \begin{cases} k & , a \leq n \leq b \\ 0 & , \text{otherwise} \end{cases}$$



## Properties of Pdf

1.  $f(n) \geq 0$
2.  $\int_{-\infty}^{\infty} f(n) dn = 1$

Ans 1 :- (i)  $f(n) = \frac{1}{b-a} = \frac{1}{10-0} = \frac{1}{10}$

$$\therefore f(n) = \begin{cases} 1/10, & 0 < x < 10 \\ 0, & \text{else} \end{cases}$$

$$\begin{aligned} \text{(ii) mean} = E(x) &= \mu = \int_{-\infty}^{\infty} n f(n) dn \\ &= \int_a^b \frac{1}{b-a} \cdot n dn \\ &= \frac{1}{2(b-a)} [n^2]_a^b \\ &= \frac{1}{2(b-a)} b^2 - a^2 \end{aligned}$$

$$\text{mean} = \frac{b+a}{2}$$

$$\therefore \text{mean} = \frac{10+0}{2} = 5$$

$$\text{variance} = E(x^2) - (E(x))^2$$

$$E(x^2) = \int_a^b \frac{1}{b-a} n^2 dn$$

$$= \left[ \frac{n^3}{3(b-a)} \right]_a^b = \frac{b^3 - a^3}{3(b-a)}$$

$$E(x) = \frac{(b-a)}{2}$$

$$E(x^2) = \frac{(b-a)(b^2+a^2+ab)}{3(b-a)}$$

$$(E(x))^2 = \frac{(b+a)^2}{4} = \frac{a^2+b^2+2ab}{4}$$

$$\text{Variance} = \frac{b^2+a^2+ab}{3} - \frac{a^2+b^2+2ab}{4}$$

$$= \frac{4b^2+4a^2+4ab-3a^2-3b^2-6ab}{12}$$

$$= \frac{b^2+a^2-2ab}{12}$$

$$= \frac{(b-a)^2}{12} = \frac{100}{12}$$

$$\text{std deviation} = \frac{10}{\sqrt{12}}$$

Ans 2: mean =  $1 = \frac{b+a}{2}$

$$b+a=2$$

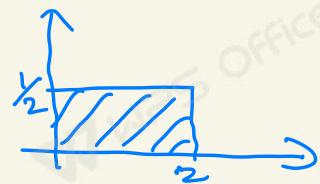
$$\text{variance} = \frac{1}{3}$$

$$\frac{(b-a)^2}{12} = \frac{1}{3}$$

$$b-a=2$$

$b=2$
$a=0$

$$P(X < 1/2) = \frac{1}{4}$$



Ans 3 :-  $f(n) = \frac{x}{4} (4-n^2)$  for  $0 \leq n \leq 2$   
= 0 otherwise

$$\Rightarrow \int n \cdot \frac{x}{4} (4-n^2) dn$$

$$\begin{aligned} E(x) &= \frac{1}{2\pi} \int_0^{2\pi} \cos(2\pi t_1 - 2\pi t_2) - \cos(2\pi t_1 + 2\pi t_2 + 2\phi) \cdot d\phi \\ &= \frac{1}{2\pi} \left[ \sin(2\pi t_1 - 2\pi t_2) - \sin(2\pi t_1 + 2\pi t_2 + 2\phi) \right]_0^{2\pi} \\ &= \frac{1}{2\pi} \end{aligned}$$

## continuous Pd

$$\mu = \frac{b+a}{2}$$

$$\text{variance} = \frac{(b-a)^2}{12}$$

$$\sigma = \frac{b-a}{\sqrt{12}}$$

$$\text{PDF} \Rightarrow f(n) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0, & \text{otherwise} \end{cases}$$

$$\text{Variance} = E(x^2) - E(x)^2$$

$$= d \int_0^d n^2 e^{-dn} dn$$

$$= d \left[ n^2 \cdot \frac{e^{-dn}}{-d} + \int 2n \cdot \frac{e^{-dn}}{\alpha} dn \right]$$

$$= d \left[ n^2 \frac{e^{-dn}}{-d} + \left[ 2n \cdot \frac{e^{-dn}}{-\alpha^2} - \int 2x \cdot \frac{e^{-dn}}{\alpha^2} dn \right] \right]$$

$$= -n^2 e^{-dn} - \frac{2n e^{-dn}}{\alpha^2} + \frac{2 e^{-dn}}{\alpha^2} \Big|_0^\infty$$

$$\frac{2}{\alpha^2}$$

## Exponential Pd

$$f(n) = \begin{cases} \alpha e^{-\alpha n}, & n \geq 0 \\ 0, & n < 0 \end{cases}$$

$$\mu = E(x) = \int_{-\infty}^{\infty} n f(n) dn$$

$$= \int_0^{\infty} n \cdot \alpha e^{-\alpha n} dn$$

$$= n \cdot \frac{\alpha e^{-\alpha n}}{-\alpha} + \int \frac{\alpha e^{-\alpha n}}{\alpha} dn \Big|_0^{\infty}$$

$$= -n e^{-\alpha n} - \frac{e^{-\alpha n}}{\alpha} \Big|_0^{\infty}$$

$$= 0 - \left[ 0 - \frac{1}{\alpha} \right]$$

$$\boxed{\mu = 1/\alpha}$$

$$\text{Variance} = \frac{2}{\alpha^2} - \frac{1}{\alpha^2} = \frac{1}{\alpha^2}$$

$$\sigma = \sqrt{\alpha}$$

Q. If the call duration is ed. with  $\alpha = 1/10$  then find

- (i) call duration exceeds 7 mins
- (ii) " between 3 & 5 mins
- (iii) " less than 8
- (iv) greater than avg of call duration.

$$\Rightarrow f(n) = \begin{cases} \frac{1}{10} e^{-\frac{n}{10}} & n \geq 0 \\ 0 & n < 0 \end{cases}$$

$$\begin{aligned} P(n > 7) &= \int_7^\infty \frac{1}{10} e^{-\frac{n}{10}} \cdot dn \\ &= \left. \frac{1}{10} \frac{e^{-n/10}}{-1/10} \right|_7^\infty \\ &= e^{7/10} \end{aligned}$$

$$\begin{aligned} (\text{i.}) \quad & \int_3^5 \frac{1}{10} e^{-\frac{n}{10}} \cdot dn \\ &= \left. -e^{-n} \right|_3^5 = -e^{-5/10} + e^{-3/10} \end{aligned}$$

$$(iii) \int_0^8 \frac{1}{10} e^{-n/10} dn$$

$$= -e^{-\frac{n}{10}} \Big|_0^8 = -e^{-\frac{8}{10}} - 1$$

$$(iv) \alpha = 1/10 \text{, } \mu = 10$$

$$P(X > 10)$$

Q Let  $Z$  be ERV - Mean = 1

$$P(Z > 2 | Z > 1)$$

$$\Rightarrow \frac{1}{\lambda} = 1 \Rightarrow \lambda = 1$$

$$f(n) = \begin{cases} 1 e^{-n}, & n \geq 0 \\ 0, & n < 0 \end{cases}$$

$$\begin{aligned} P(Z > 2) &= \int_2^\infty e^{-n} \cdot dn \\ &= -e^{-n} \Big|_2^\infty \\ &= +e^{-2} \end{aligned}$$

$$P(Z > 1) = \int_1^{\infty} e^{-n} dn$$

$$= -e^{-n} \Big|_1^{\infty} = e^{-1}$$

$$P(Z > z_2 | Z > 1) = \frac{e^{-z_2}}{e^{-1}} = e^{-z_2}$$

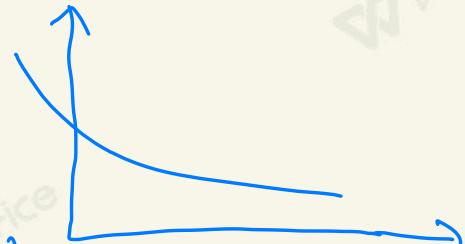
Q. Let  $X_1$  &  $X_2$  be two independent ERV. with mean = 0.5 & 0.25 respectively then  $y = \min(X_1, X_2)$ . What

$$\Rightarrow \alpha_1 = \frac{1}{0.5} = 2$$

$$\alpha_2 = \frac{1}{0.25} = 4$$

$$f(x_1) = \begin{cases} 2e^{-2x} & \\ \end{cases}$$

$$f(x_2) = \begin{cases} 4e^{-4x} & \\ \end{cases}$$



$$= \min(2e^{-2x}, 4e^{-4x})$$

=

$$\text{Var}(x+y) = \text{Var}(x) + \text{Var}(y) + 2\text{Cov}(x, y)$$

if  $x$  &  $y$  are independent variable,

$$\text{Var}(x+y) = \text{Var}(x) + \text{Var}(y)$$

$$\text{correlation coeff} = \rho = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x) \cdot \text{Var}(y)}}$$

Q. consider two boxes box1 & box2. Box 1 contains 4 red and 6 black balls whereas Box2 contains 5 red and 5 black balls. Now, a coin is tossed, if head occurs then one ball is randomly drawn from box1. whereas if tail occurs then one ball from Box2.

- (i) Find the probability of getting a red ball.
- (ii) if the ball obtained is red, what is the prob that it comes from box1.

$$\Rightarrow P(B_1) = 1/2 = P(B_2)$$

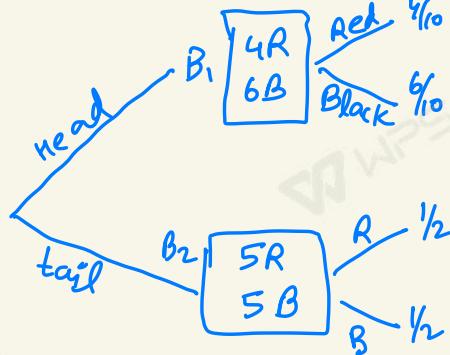
$$P(\text{Red}) = P(\text{Red} \cap B_1) + P(\text{Red} \cap B_2)$$

$$P(\text{Red}) = P(A) \cdot P(E|A) + P(B) \cdot P(E|B)$$

$$= \frac{1}{2} \times \frac{4}{10} + \frac{1}{2} \times \frac{1}{2}$$

$$= 9/20$$

$$(ii) \quad \frac{\text{fav. path}}{\text{total path}} = P(B_1 | \text{Red}) = \frac{P(B_1) \cap P(\text{Red})}{P(\text{Red})} = \frac{\frac{1}{2} \times \frac{4}{10}}{\frac{9}{20}} = \frac{4}{9}$$



## Types of Events :-

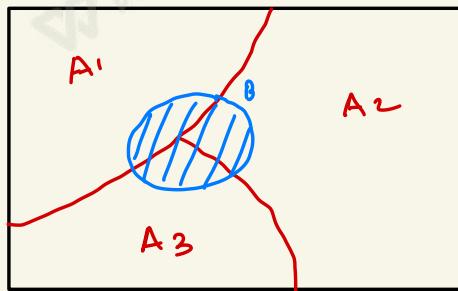
- 1) Mutually dependent
- 2) Mutually independent ( $P(A \cap B) = P(A) \cdot P(B)$ )
- 3) Mutually Exclusive (either head or tail)
- 4) Collectively exhaustive ( $P(A) + P(B) + P(C) = 1$ )

$$J = M \cdot C$$

Joint prob | marginal prob  
 (independent) | conditional  
 (dependent)

$$P(A \cap B) = P(A) \cdot P(B|A)$$

→ If  $P(A \cap B) = 0$ , then A & B are mutually exclusive.



$$\begin{aligned}
 P(B) &= P(A_1 \cap B) + P(A_2 \cap B) + P(A_3 \cap B) \\
 &= P(A_1) \cdot P(B|A_1) + P(A_2) \cdot P(B|A_2) + P(A_3) \cdot P(B|A_3)
 \end{aligned}$$

$$P(B) = \sum_{i=1}^n P(A_i) \cdot P(B|A_i)$$

Total probability theorem

## Baye's theorem

$$P(A_2|B) = \frac{P(A_2 \cap B)}{P(B)} = \frac{P(A_2) \cdot P(B|A_2)}{P(B)}$$

$$P(A_2|B) = \frac{P(A_2) \cdot P(B|A_2)}{P(B)}$$

Q. In a bolt factory, machines A, B & C manufacture 25%, 35% & 40% respectively of the total output. There is a chance of 5%, 9% and 2% resp are defective. A bolt is drawn at random.

- (i) Prob that it is a defective bolt 0.0345  
(ii) Prob that defective is drawn from A. 0.3623

$$\begin{aligned} \therefore P(\text{def}) &= P(\text{def} \cap A) + P(\text{def} \cap B) + P(\text{def} \cap C) \\ &= \frac{25}{100} \times \frac{5}{100} + \frac{35}{100} \times \frac{9}{100} + \frac{40}{100} \times \frac{2}{100} \\ &= \frac{125 + 315 + 80}{10000} \\ &= \frac{520}{10000} = \end{aligned}$$

- Q. Consider the following corpus of 4 sentences.
- (S) those friends Aman, akbar and anthony are reading book. (S)
- (S) aman is reading malgudi days (S)
- (S) akbar is reading a detective book (S)
- (S) anthony is reading a book by nk narayan (S)
- Assume a bigram language model. Calculate  
 $p(S) \text{ aman is reading a book } (S))$

Q. There are 3 bags  $B_1, B_2$  &  $B_3$ . The bag  $B_1$  contains 5 green & 5 red balls, bag  $B_2$  contains 3 red & 5 green balls and bag  $B_3$  contains 5 red & 3 green balls. Bags  $B_1, B_2$  &  $B_3$  have probabilities  $\frac{3}{10}, \frac{3}{10}$  &  $\frac{4}{10}$  respectively if being chosen. A bag is selected at random and a ball is chosen at random from bag.

- (A) Find the probability that the chosen ball is green, given that the selected bag is  $B_3$ .
- (B) " " " " " is green.
- (C) Find the probability that the selected bag is  $B_3$ , given that the chosen ball is green.
- (D) Find the prob. that the selected Bag is  $B_3$  & given that the chosen ball is green.

Ans 2 :-

(B)  $P(\text{green}) = \frac{3}{10} \times \frac{5}{10} + \frac{3}{10} \times \frac{5}{8} + \frac{4}{10} \times \frac{3}{8}$

$$= 0.15 + \frac{15}{80} + \frac{12}{80}$$

$$= 0.15 + 0.1875 + 0.15$$

$$= 0.4875$$

(A) =

$$\frac{4/10 \times 3/8}{4/10} = 3/8$$

(C)  $P = \frac{4/10}{0.4875} = 0.8205$

(D)  $P = \frac{4}{10} \times \frac{0.4875}{0.195}$

$$P(B_3 \cap \text{Green}) = P(B_3) \cdot P(\text{green}|B_3) = 3/10$$

Ans 1 :-

$$P = \frac{3+6+8+4+4}{24} = \frac{25}{24}$$

Q. Consider the following corpus of 3 sentences what is the total count of unique bigrams for which the likelihood will be estimated assume we do not perform any preprocessing. consider the beginning & end of token as <S> & </S>

→ <S> Julia is visiting the museum </S>  
<S> Julia, grover & natasha are friends </S>  
<S> zoe & natasha will meet julia in  
the museum </S>

(a) 23

✓ (b) 20

(c) 16

(d) 18

## Function words vs content words

stop words  
or  
closed class words  
a, an, the, to, is, of, ---  
prepositions - in, on, of, by  
determiners - a, an, the  
pronouns - I, he, his, him, me, the

↓ info / topic  
open class words  
↳ we keep on adding  
new words.  
noun  
verb  
adjective  
adverbs

## Type Token Ratio (TTR)

$$TTR = \frac{\text{types}}{\text{tokens}} = \frac{\text{no. of unique words}}{\text{total words}}$$

Ex:- will will will

$$TTR = 1/3$$

\* If TTR is high, then new words will be found more.

Q. In a corpus, it was found that the word with rank 4<sup>th</sup> has a frequency of 600. What can be the best guess for the rank of a word with freq. 300.  $\Rightarrow 2600 \times 4 = 300 \times n$   $n = 8$

Q. In the sentence, "The only thing we have to fear is the fear itself", Find the TTR.  $\Rightarrow \frac{9}{11}$

Q. Let the rank of 2 words  $w_1$  &  $w_2$  in a corpus be 1600 and 400 respectively. Let  $m_1$  &  $m_2$  represent the no. of meanings of  $w_1$  and  $w_2$  respectively. The ratio of  $m_1 : m_2$  would tentatively —

$$m_1 \times \sqrt{1600} = m_2 \times \sqrt{400} \quad 1/2$$

Q. Which are true :-

T (a) Ambiguity can appear in Tokenization steps.  
F (b) Ambiguity will not appear in sentence segmentation step.

T (c) Function used is generally more frequent in a text than any any content word.

T (d) Output of lemmatization are always real words

### Zipf's Law

$$f \propto \frac{1}{r}$$

f: freq of word  
r: rank of word  
in dec ord

$$\therefore f \cdot r = k = \text{constant}$$

$$P_r (\text{prob of word or rank } r) = \frac{f}{N} = \frac{\frac{1}{r}}{\sum_i \frac{1}{r_i}}$$

$$= \frac{1}{r}$$

$$m \propto \sqrt{f}$$

m: no. of meanings

$$m \propto \frac{1}{\sqrt{f}}$$

$\propto$ : rank

$$f \propto \frac{1}{l}$$

l: length of a word

## Heap's Law

$$|U| = KN^\beta$$

$|U|$  = size of vocabulary  
 $N$  = number of tokens

Q. what is the size of unique words in a document where total no. of words is 12000,  $K=3.71$  &  $\beta=0.69$ .

∴

$$|U| = 3.71 \times (12000)^{0.69}$$
$$= 2421$$

Q. If the first corpus has  $TTR = 0.085$  & second corpus has  $TTR = 0.78$ . Which of the following are F.

F (i) 1<sup>st</sup> corpus has more tendency to have unique words.

T (ii) 2<sup>nd</sup> .. .. .. .. ..

F (iii) TTR values can have sometimes  $> 1$ .

T (iv) ↑ TTR indicated ↑ degree of lexical variance & vice versa.

1) Ambiguity in Lexicography

2) Issues in Tokenization

10-12 slides

# Conversion of text into vectors

(vectorization)

for  
feature  
extraction  
or  
Information  
retrieval

1. Bow  $\rightarrow$  Bag of words model

2. Binary Bow

3. TF-IDF (Term frequency-Inverse doc freq)

ML  $\rightarrow$  label encoding & OHE (one hot encoding)

Corpus

$S_1$ : He is an awesome boy and an awesome dancer too.

$S_2$ : She is also an awesome girl

$S_3$ : Both the girl and the boy are awesome.

$S_4$ : It was a good movie with awesome acting.

After removal of stopwatch

Vocabulary		freq
awesome	$f_1$	5
boy	$f_2$	2
dance	$f_4$	1
girl	$f_3$	2
good	$f_5$	1
movie	$f_6$	1
acting	$f_7$	1

Text pre-process.

① Stopwords removal

② Lowering of the case

③ Stemming & lemmatization

④ Removal of punctuation sym.

⑤ handling of negation

vectors for 1st sentence

	f <sub>1</sub>	f <sub>2</sub>	f <sub>3</sub>	f <sub>4</sub>	f <sub>5</sub>	f <sub>6</sub>	f <sub>7</sub>
awesome	1	0	1	0	1	0	0
[ 2 ]	1	0	1	0	1	0	0 ]

Vectors → fixed size

- Q. S<sub>1</sub>: The pizza was good  
 S<sub>2</sub>: The pizza was not good

→

	good	the	pizza	was	not
S <sub>1</sub> :	1	1	1	1	0 ]
S <sub>2</sub> :	1	1	1	1	1 ]

TF(t, d) = total no. of times 't' present in doc. d  
total no. of tokens in the given docd

(Term freq.)

↓

$$IDF(t) = \log \left( \frac{\text{Total no. of documents}}{\text{no. of documents containing the term t}} \right)$$

(Inverse document freq.)

TF-IDF = TF(t, d) × IDF(t)  
 ↓ score

- Q. d<sub>1</sub>: He is a good boy  
 d<sub>2</sub>: She is a good girl  
 d<sub>3</sub>: Both the boy and the girl are good  
 Calculate TF-IDF score for the above corpus.

→ d<sub>1</sub>: good boy  
 d<sub>2</sub>: good girl  
 d<sub>3</sub>: boy girl good

	$d_1$	$d_2$	$d_3$	
good:	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{3}$	(TF( $t_{sd}$ ))
boy:	$\frac{1}{2}$	0	$\frac{1}{3}$	
girl:	0	$\frac{1}{2}$	$\frac{1}{3}$	

$$\text{good: } \log \frac{3}{3} = 0$$

$$\text{boy: } \log \frac{3}{2}$$

$$\text{girl: } \log \frac{3}{2}$$

$$(\text{IDF}(t))$$

	<u>good</u>	<u>boy</u>	<u>girl</u>	
$d_1:$	0	$\frac{1}{2} \times \log \frac{3}{2}$	0	(TF-IDF)
$d_2:$	0	0	$\frac{1}{2} \log \frac{3}{2}$	$(d_i \times \text{IDF}_i)$
$d_3:$	0	$\frac{1}{3} \log \frac{3}{2}$	$\frac{1}{3} \log \frac{3}{2}$	

- Q. From a website we got 3 reviews of a movie:
- R1: This movie is very scary and long.
  - R2: This movie is not scary and is slow.
  - R3: This movie is spooky and good.

	$d_1$	$d_2$	$d_3$	$\text{IDF}$
This:	$\frac{1}{7}$	$\frac{1}{8}$	$\frac{1}{6}$	$\log \frac{3/3}{3} = 0$
movie:	$\frac{1}{7}$	$\frac{1}{8}$	$\frac{1}{6}$	$\log \frac{3/3}{3} = 0$
is:	$\frac{1}{7}$	$\frac{2}{8}$	$\frac{1}{6}$	$\log \frac{3/3}{3} = 0$
very:	$\frac{1}{7}$	0	0	$\log \frac{3/1}{1}$
scary:	$\frac{1}{7}$	$\frac{1}{8}$	0	$\log \frac{3/2}{2}$
and:	$\frac{1}{7}$	$\frac{1}{8}$	$\frac{1}{6}$	$\log \frac{3/3}{3} = 0$
long:	$\frac{1}{7}$	0	0	$\log \frac{3/1}{1}$
not:	0	$\frac{1}{8}$	0	$\log \frac{3/1}{1}$
slow:	0	$\frac{1}{8}$	0	$\log \frac{3/1}{1}$
spooky:	0	0	$\frac{1}{6}$	$\log \frac{3/1}{1}$
good:	0	0	$\frac{1}{6}$	$\log \frac{3/1}{1}$

R <sub>i</sub> :	<u>is</u>	<u>This</u>	<u>movie</u>	<u>and</u>	<u>Scary</u>	<u>very</u>	<u>long</u>	<u>spooky</u>	<u>slow</u>	<u>not good</u>
R <sub>j</sub> :	0	0	0	0	$\frac{1}{7} \log_2 3$	$\frac{1}{7} \log_2 3$	$\frac{1}{7} \log_2 3$	0	0	$\frac{1}{8} \log_2 3$
R <sub>j</sub> :	0	0	0	0	0	0	0	$\frac{1}{6} \log_2 3$	0	$\frac{1}{6} \log_2 3$

## Spelling correction - Edit Distance

I am writing an email on behaf of KIIT.  
incorrect

min. distance

↓  
min no. of operations

↳ insertion

↳ deletion

↳ substitution

↳ correct one

• behalf

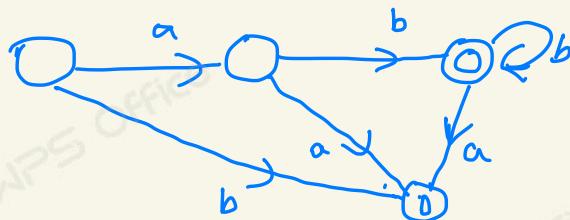
• behave

• behaviour

} Levenshtein  
distance

daqb}

abb\*--



N	9	8	9	10	11	12	11	10	9	8
O	8	7	8	9	10	11	10	9	8	9
I	7	6	7	8	9	10	9	8	9	10
T	6	5	6	7	8	9	10	11	10	9
H	5	4	5	6	7	8	9	10	9	8
E	4	3	4	5	6	7	8	7	9	8
T	3	4	5	6	7	8	7	8	7	7
N	2	3	4	5	6	7	8	7	6	7
I	1	2	3	4	5	6	7	6	7	8
#	0	1	2	3	4	5	6	7	8	9
	#	E	X	E	C	U	T	I	O	N

$$D(i,j) = \min \begin{cases} D(i-1, j) + 1 \\ D(i, j-1) + 1 \\ D(i-1, j-1) + n \end{cases}$$

$$\begin{aligned} n &= 2 && \text{if } x(i) \neq y(j) \\ &= 0 && \text{if } x(i) = y(j) \end{aligned}$$

Q. "The boy put the keys on the table". Assign the POS tag for each word.

⇒ The → det  
 boy → noun  
 put → verb  
 the → det  
 keys → noun  
 on → pre  
 the → det  
 table → noun

Q. "The grand jury commented on a number of other topics".

⇒ The → DT (determiner)  
 grand → JJ (adjective)  
 jury → NN (noun)  
 commented → VBD (verb)  
 on → pre - IN  
 a → det - DT  
 number → noun - NN  
 of → Prep - IN  
 other → JJ (adjective)  
 topics → NNS (singular noun)

Q. "I need a flight from Atlanta".

→ I → PRP (personal pronoun)  
need → VB (verb)  
a → DT (determiner)  
flight → NN  
from → IN  
Atlanta → NNP

$$P(\text{VB} \mid \text{to}) \times P(\text{NR} \mid \text{VB}) \times P(\text{race} \mid \text{VB})$$

state      STP      emission prob.  
transition prob.

Q. (1) Sun rises in the east.  $P=1, I=0$

(2) India got Independence in 1947.  $P=1, I=0$

(3) You will get 5 emails in next one year.  $P=0.04 \quad I\uparrow$

(4) The prime minister will come to your home tomorrow.  $P=0.01 \quad I\uparrow$

(5) It will snow in Delhi in June.  $P=0.002 \quad I\uparrow$

(6) You will not get a holiday next Sunday.  $P=0.04 \quad I\uparrow$

(7) The dog will bark.  $P \uparrow \quad I\downarrow$

$$P \propto \frac{1}{I}$$

I: Information content

$$P(n_i) \propto \frac{1}{I(n_i)}$$

$$P(n_i) = f\left(\frac{1}{I(n_i)}\right)$$

$$P(n_i) = \log \frac{1}{I(n_i)}$$

$$= \log (I(n_i))^{-1}$$

$$P(n_i) = -\log I(n_i)$$

$$P(n_i) P(y_i) = I(n_i) + I(y_i)$$

$$\text{Entropy } H = \sum_{i=1}^n P(n_i) I(n_i)$$

$$= \sum_{i=1}^n P(n_i) \cdot \log \frac{1}{P(n_i)}$$

$$H = - \sum_{i=1}^n P(n_i) \log P(n_i)$$

$$\sum P(n_i) \log \frac{1}{P(n_i)}$$

$$L_{\min} = H$$

$$L_{avg} = \sum p_i L_i$$

$$P(x, y) = P(x) \cdot P(y|x)$$

↓  
Independent

$$\therefore H(x, y) = H(x) - H(y|x)$$

chain rule of  
entropy

↓  
Independent