

Perfect Computer Engineer

WARNING

It was found that Students purchased these notes and contaminate them by mixing wrong information in the answers, deleting the answers or adding new answers with completely wrong information and sharing it so that the other Student who is getting these notes and who also happens to be a competitor of this student score very less exam marks. Perfect Computer Engineer YouTube channel doesn't have any problem with sharing the notes. (Channels main aim is to help students get high exam marks) but there are high chances of you getting notes with highly contaminated information if you are getting them from any social media sites, telegram, WhatsApp, Discord, or from someone else.

If this is the case Perfect Computer YouTube channel doesn't take any guarantee or responsibility of the authenticity and originality of the information in these notes.

There are only two sources of purchasing these Notes:

1. Directly from the Seller: Ojas Deshpande. Instead: planetofes
2. From the YouTube Perfect Computer Engineer (check the description of any video related to the topic in this subject)

Thank You So much...

Ojas Deshpande

Natural Language Processing

* Introduction to Natural Language Processing

Topic

Page No.

1. Introduction to Natural Language Processing

1

2. Steps in Natural Language Processing

2-4

3. Stages/Levels in Natural Language Processing

5-6

4. Ambiguity in Natural Language Processing

7-8

5. Applications of NLP

9-10

* Word Level Analysis

1. Morphological Parsing & Morphological analysis

11-13

2. Inflectional & Derivational Morphology.

14-16

3. Finite State Automata
[Numericals]

17-20

Page No.

21

4. Inflectional Vs Derivational Morphology

5. Language Model

22-24

6. N-Gram Model

24-25

7. Numerical on bigram

26

8. Numerical on predict the next word

27-28

* Syntactic Analysis

29

1. Part of Speech tagging

30

2. Why Part of Speech tagging

31

3. Why is tagging hard

31

4. Applications of PoS Tagging

32

5. Rule Based POSTagger

34

6. Stochastic PoS Tagger

36-37

7. Transformation based Tagger

38 - 39

8. Advantages and Disadvantages of Transformation based tagging

39-40

9. Multiple Tags, Multiple words
Unknown words. 41-43

10. Context Free Grammar 44-45

11. Numerals to CFG rules 46

12. Top down and bottom
up parsing Numericals. 47-48

13. Top Down VS Bottom
up Parsing 49

* 15 Semantic Analysis

1. Introduction 50-51

2. Applications of semantic
Analysis 52-53

3. Wordnet 54

4. Structure of wordnet 55

5. Applications of wordnet 56-57

6. Word sense Disambiguation 58 - 62

7. 30 Homonymy, Polysemy,
Hyperonymy, Hyponymy,
Meronymy, Synonymy,
Antonymy.

* Pragmatics

1. Introduction & Discourse. 69 - 72

* Applications

- 1 Machine Translation 73 - 75

2. Information Retrieval 76 - 77

3. Question Answer System 78 - 80

4. Text Categorization 81 - 84

5. Text summarization 85 - 87

* Book Recommendations

* Special thanks for the betterment of the notes in terms of knowledge, Structure, content etc

1. Mrs. Rajeshree

2. Internet (Random sizes)

3. TechKnowledge

4. Mr. Varun

5. LMT

6. YouTube (Random channels)

7. Dr. Sharvari etc...

Module 1 : Introduction to Natural Language Processing

Q. 1. Introduction to Natural Language Processing ?

⇒ The process of computer analysis of input provided by human language (natural language), and conversion of this input into a useful form of representation is Natural Language Processing.

10. The field of Natural Language Processing is primarily concerned with getting computers to perform useful and interesting tasks with human languages.

15. The field of Natural Language Processing is secondarily concerned with helping us come to a better understanding of human language.

There are few major components in Natural Language Processing.

1. Natural Language Understanding

- Mapping the given input in the natural language into a useful representation.

- There are different levels of analysis required here :

Morphological Analysis, Syntactic Analysis, Semantic analysis, discourse analysis etc.

2. Natural Language Generation

• Producing output in the natural language from some internal representation.

- Different level of Synthesis required : deep planning (what to say), syntactic generation.

5 Goals of Natural Language Processing

- Design, implement and test systems that process natural languages (Ex: Marathi) for practical applications.

Q2 10 Steps in Natural Language Processing

1] Tokenization

- o Process of cutting big sentence into small tokens.
- o Example :
[Dipesh] [is] [an] [abnormal] [human]

2] Stemming

- o Normalizing words into their base or root forms

o Example : Knows Known Knowing

Know
→ root form

Disadvantage : Sometimes root word is not sensible.

Example : Date Dated Dating

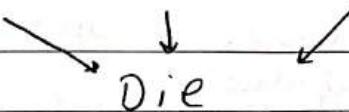
Dat

3] Lemmatization:

Groups together different inflected forms of a word called Lemmatization. This is similar to stemming but here the output is always correct.

Example:

Die Died Dead



4. POS Tags:

This stands for Part of Speech Tags. It indicates how a word functions in meaning as well as grammatically within a sentence.

Manan killed a bat and ate it

↓ ↓ ↓ ↓ ↓ ↓ ↓

Noun Verb det Noun conj verb pronoun

unction
ion

Risadvantages: One word can have multiple part of speech, we use Name entity recognition for this problem.

5. Name Entity Recognition:

It is a process of putting Name entity on the words: person, organization, company, location etc.

e.g. google it

company

6] Chunks

picking individual pieces of information and grouping them into bigger pieces.

5

Anurag ate the black cat

Noun verb Det Noun Noun

i , |] []]

Noun Verb Noun Phrase
phrase phrase

↓ ↘
chunk

10

A Chunk helps in getting insightful and meaningful info from the text.

15

Q3 Stages / Levels of Natural Language Processing

=> Morphology : concerns the way words are built up from smaller meaning bearing units

example : " truthfulness"

Syntax : concerns how words are put together to form correct sentences and what structural role each word has.

Example : " The dog ate my homework .

Not syntactically correct .

Semantics : concerns what words mean and how these meanings combine in sentences to form sentence meanings .

Eg : industrial plant / living organism
Doesn't make a proper meaning .

Pragmatics : concerns how sentences are used in different situations ... and how it affects the interpretation of the sentence .

Sentence might be correct but it might have semantic ambiguity - which means that sentence has two meanings .

Eg : Dipesh loves his girlfriend and Mamam does too

Page : 6

Date :

5 Discourse: concerns how the immediately preceding sentences affect the interpretation of the next sentence.

Q4. Write a short note on Ambiguity

Ambiguity in Natural Language Processing can be referred to as the ability of being understood in more than one way.

Natural language is very ambiguous.

Natural Language Processing has following ambiguities.

1. Lexical Ambiguity:

Ambiguity of a single word is called Lexical Ambiguity.

Example: I can play cricket

give me that can

both "can" have different meanings.

2. Syntactic Ambiguity:

This kind of ambiguity occurs when sentence is parsed in different ways.

Ex: Abid saw the man with the binoculars

- Abid saw the man carrying binoculars

- Abid saw the man through the binoculars.

3. Semantic Ambiguity :

This occurs when actual or erad meaning of the phrase themselves can be mismatched or misinterpreted when even after syntax and meaning of individual word have been resolved.

Eg : Manam loves his cat and divesh does too.

4. Anaphoric Ambiguity :

This kind of ambiguity arises due to use of anaphora entities in discourse

Anaphora : When a noun is replaced by a pronoun and causes confusion

Eg : The house is on a long street. It is very dirty.

5. Pragmatic Ambiguity : It occurs when a sentence gives its multiple interpretation or it is not specific

Eg : I love you too

Note : Please Refer my video lectures, if not understood

Q85 Applications of Natural Language Processing

=> 1. Question Answering - by computer.

Question Answering (QA) system is a task of automatically answering to the questions asked in the natural language using either a pre-structured database or a collection of natural language documents.

It presents only the requested information instead of searching full documents like search engine.

The basic idea behind this is user just have to ask the question and the system will retrieve the most appropriate and correct answer.

Eg. Q: who is the father of the Nation [India]
A. Mahatma Gandhi

2. Chatbots:

Intelligent chatbots are offering personalised assistance to the customers already. Analysts predict that the use of chatbots will grow 5 times year on year.

3. Managing the Advertisement Funnel

NLP is a great source for intelligent targeting and placement of advertisements in the right place at the right time and for the right audience.

4. Market Intelligence

NLP gives exhaustive insights into employment changes and status of the market, tender delays, and closings which help in extracting information from large repositories.

5. Text Summarization:

10

It refers to the technique of shortening long pieces of text.

15

Summarization can be mainly classified into extractive and abstractive.

→ Extractive summarization → Extracting few sentences

→ Abstractive summarization → It builds an internal semantic representation and then uses natural language generation techniques to create a summary.

25

In both of the cases NLP is used.

30

Module 2 : Word Level Analysis

Q. Write a short note on morphological parsing and Morphology Analysis.

=> Morphological parsing is the task of recognizing the morphemes inside a word.

- Morphemes are the minimal meaning-bearing unit in a language.

10

- Example : Mangoes

Mango es → ①

Here there are two Morphemes

15 Morphemes can be Stems (Root word) or an Affix

- Now this Affix is divided into three parts
An affix can be prefix (eg reform) or
20 Suffix (eg loved) or infix (passersby).

- So here Mango is a Stem and es is a suffix because it is attached after the main word

25 Following are the requirements for building a morphological Parser

1. Lexicon : It includes the list of stems and affixes along with the basic information about them.

30 eg Stem is a noun stem or a verb stem.

Morphotactics :

Morphotactics has a set of rules by the help of which it decides the ordering of words.

example : USE able ness

 \ / /
 USEability ness ✓

able use ness
 \ / /
 ableness ness ✗

Orthographic rules :

These spelling rules are used to model the changes occurring in a word.

example : lady + s = ladys ✗

lady + s = ladies ✓

- The study of formation of words is called morphology.
- Some words are self sufficient they have their own meaning example - camera, pen

- + Some words if divided has there own meanings example : Showcase = Show & Case both words have there own meanings.
- Some words if combine don't have any meaning but if they are combined with a word it becomes a meaningful word.
example: ing has no meaning but if combined with love its loving which have there meaning.
- So basically there are different words which if used in a right way we can get a meaningful word.
- So analysis is studying in detail and Analysis of Morphology is Morphology analysis.

Q1 write a short note on Inflectional and Derivational morphology.

=> Before Inflectional and derivational morphology we need to understand what are morphemes. Morpheme is a word or a part of a word that has meaning and a morpheme cannot be divided further into meaningful units.

example of morpheme : cat

If we try to divide morpheme more it will be a meaningless result.

There are two types of morphemes

① Free Morphemes ② Bound morphemes

① Free Morpheme : free morpheme is a morpheme which has its own meaning or it has its complete meaning example : fan, camera etc

- Free Morphemes are of two types lexical morphemes and grammatical morphemes

- Lexical morphemes are the picture words they are noun, adjective, verbs, adverbs

example : black, yellow, chair

every year new and new lexical morphemes are added in a language.

- Grammatical morphemes are grammar words which are limited in each and every language.

- These morphemes don't change frequently like Lexical morphemes they are preposition, conjunctions, etc

② Bound Morphemes : These morphemes are of two types Inflectional and Derivational. But before going further we must know what is bound morphemes.

Bound morphemes are those morphemes whose meaning is not complete in themselves. And that is the reason why they depend on the free morphemes for meaning.

Now Affixes are bound morphemes and Affixes are of three types prefix, suffix, prefix
prefix → Because

Suffix → loveable

Infix → passersby

Inflectional morpheme : Inflectional morpheme is one which when attached to a root word doesn't change its class

BooR + S = Books
Noun Noun

Inflectional morphemes are Infixes and Suffixes and can't be prefix

Derivational Morphemes:

Derivational Morphemes are ones which when added to a word changes its class.

Teach + er = Teacher
Verb Noun

10 Derivational morphemes are of two types class changing which was the above one and class maintaining which when added to word changes the word but can't change the class.

20

25

30

Q3 Design a finite state automata (FSA) for $b \circ a + !$

$\Rightarrow Q$: finite set of states

q_0, q_1, q_2, q_3, q_4

Σ : set of input alphabets
 $\{a, b, !\}$

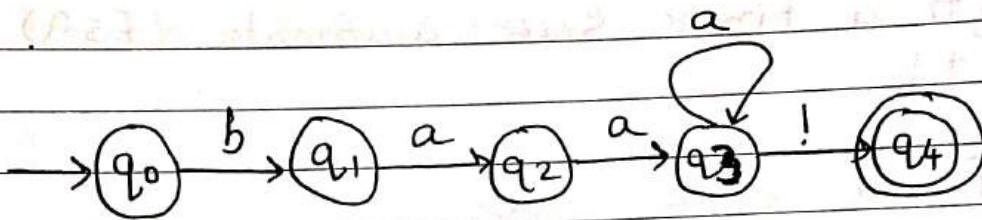
q_0 : Start State.

F: Set of final states
 $F \subseteq Q$

$S(q, i)$: the transition function or transition matrix between states. Given a state $q \in Q$ and input symbol $i \in \Sigma$, $S(q, i)$ returns a new state $q' \in Q$. S is thus a relation from $Q \times \Sigma \rightarrow Q$;

Transition table

States \ input	a	b	!
States	\emptyset	q_1	\emptyset
q_0	q_2	\emptyset	\emptyset
q_1	q_3	\emptyset	\emptyset
q_2	q_3	\emptyset	\emptyset
q_3	q_3	\emptyset	q_4
q_4	\emptyset	\emptyset	\emptyset



Q4 Design a finite State Automata for divisibility by 5 tester for binary number.

$Q \Rightarrow$ finite set of state

$q_0 \Rightarrow$ Start state ("Might or might not include your choice")

$q_0 \Rightarrow$ rem 0 State

$q_1 \Rightarrow$ rem 1 State

$q_2 \Rightarrow$ rem 2 State

$q_3 \Rightarrow$ rem 3 State

$q_4 \Rightarrow$ rem 4 State

rem = remainder.

The question can also be written as design a FSA to check whether given binary no is divisible by 5 or not

$\Sigma \Rightarrow$ set of input alphabets.

$$= \{0, 1\}$$

$S =$ Transition function

$$S : Q \times \Sigma \rightarrow Q$$

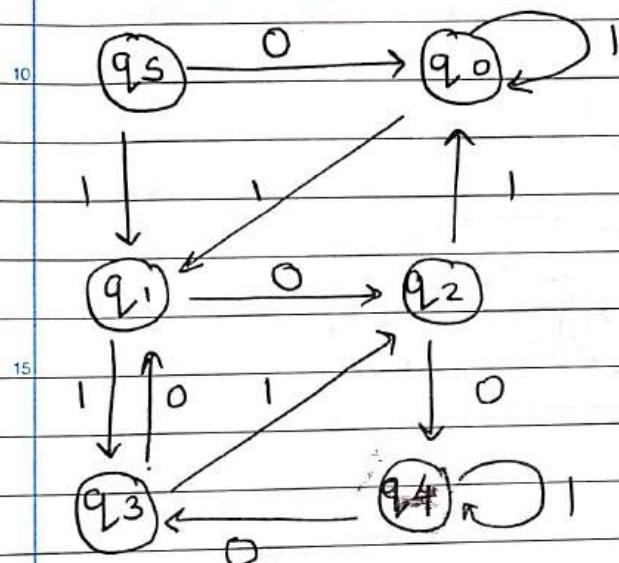
$q_0 =$ Start state = q_S

F = final state

$$q_0, F \subseteq Q$$

[or you can write like previous Numerical]

$s \setminus I$	0	1
q_s	q_0	q_1
*	q_0	q_1
q_1	q_2	q_3
q_2	q_4	q_0
q_3	q_1	q_2
q_4	q_3	q_4



Q5. Design a DFA of a string that should end with 100

$$\Rightarrow M = \{ Q, \Sigma, S, q_0, F \}$$

q_0 = initial state

q_1 = String ending with 1

q_2 = String ending with 10

q_3 = String ending with 100

$$\Sigma = \{ 0, 1 \}$$

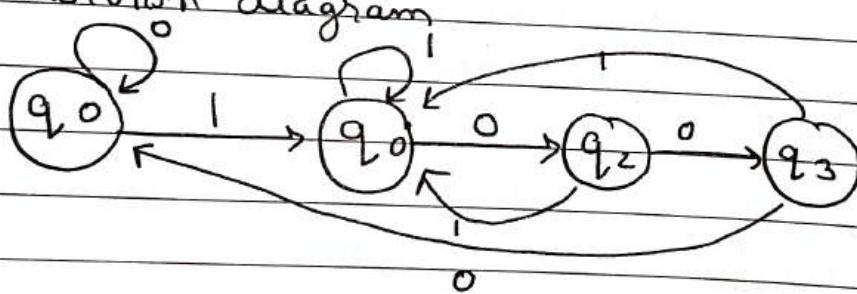
$$q_0 = \{ q_0 \}$$

F = final state = q_3

Transition Table

	0	1
5	q_0	q_0
10	q_1	q_2
	q_2	q_1
	q_3	q_1

Transition diagram



Q6 Differentiate between Inflectional & Derivational morphology.

5. Inflectional Morphology

1. It is a morphological process that adapts existing words so that they function efficiently in sentences without changing POS of base morpheme.

Derivational Morphology

It is concerned with the way morphemes are connected to existing lexical forms as affixes.

2. Regular : It is more regular

It is very less regular

3. Use : Can only be Suffix or infix and not prefix

can be both prefix & suffix

4. Change in : Never changes part of the grammatical category or POS

It can change the grammatical category or POS

5. Example: Cat + S = Cats
Noun Noun

danger + ous = dangerous
Noun Adjective.

Q7. Write a short note on language model

=> The goal of a language model is to assign probability to a sentence.

With this it also decides which sentence is more accurate at the moment.

Example : She is a tall girl is more accurate than She is a long girl

Statistical Language Modelling , Or Language Modelling is the development of probabilistic models that are able to predict the next word in the sequence given the words that precede it . It is a probability distribution over sequences of words .

Given such a sequence , say of length m , it assigns a probability $P(w_1, \dots, w_m)$ to the whole sequence .

The goal of probabilistic language modelling is to calculate the probability of a sentence of sequence of words : $P(w) = P(w_1, w_2, w_3 \dots w_n)$ and can be used to find the probability of the next word in the sequence :

$$P(w_5 | w_1, w_2, w_3, w_4)$$

A model that computes either of these is called language model

Method of calculating Probability:
conditional probability:

let A and B be two events with $P(B) \neq 0$, the conditional probability of A given B is:

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

$$P(x_1, x_2, \dots, x_n) = P(x_1) P(x_2|x_1) \dots P(x_n|x_1, \dots, x_{n-1})$$

for example: $P(\text{"its water is so transparent"}) =$
 $P(\text{its}) * P(\text{water}|\text{its}) * P(\text{is}|\text{its water}) * P(\text{so}|\text{its water is}) * P(\text{transparent}|\text{its water is so})$

we can estimate this by simply counting and dividing the results.

$$P(\text{transparent} | \text{its water is so}) = \frac{\text{count(its water is so transparent)}}{\text{count(its water is so)}}$$

Markov Property : It says that the probability of the next word can be estimated given only the previous K number of words.

for example, if $K=1$:

$$P(\text{transparent} | \text{its water is so}) \approx P(\text{transparent} | \text{so})$$

or if $K=2$:

$$P(\text{transparent} | \text{its water is so}) \approx P(\text{transparent} | \text{is so})$$

~~General Assumption, the probability of next word depends on previous words.~~

general equation for the Markov Assumption,
 $k = i$:

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

Q8. Write a short note on N-Gram model.

=> Note: This answer is in continuation with
 10 the previous answer. You have to decide the
 length depending the marks allotted to question.

- The Simplest case of markov model is a unigram model, In this model we simply estimate the probability of the whole sequence of words by the product of probabilities of individual words - unigrams.
 15 and if we generated sentences by randomly picking words, it would be

20 Sixth, the, rupees, abduction

It would be just a random sequence of words

$$P(w_1, w_2, \dots, w_n) \approx \prod_i P(w_i)$$

- Slightly more intelligent is the bigram model where we condition on the single previous word ..

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

5 We can extend this to trigrams, 4-grams, 5-grams.

But in general this is an insufficient model
of language.

o because language has long distance dependencies
example:

10 "The computer which I had just put into the
machine room on the fifth floor crashed."

And if we say predict the next word after floor
So g_4 is very unlikely to predict crash but
if we compare with or bring the main subject
"computer" in the picture then we are more
likely to guess crashed or predict crash as a
next word.

Q9 corpus:

$\langle S \rangle I \text{ am a human } \langle /S \rangle$

$\langle S \rangle I \text{ am not a stone } \langle /S \rangle$

$\langle S \rangle I \text{ I live in Mumbai } \langle /S \rangle$

check the probability of $\langle S \rangle I \text{ I am not } \langle /S \rangle$
using bigram

$$\Rightarrow P(I \text{ I am not})$$

$$= P(I | \langle S \rangle) P(I | I) P(\text{am} | I) P(\text{not} | \text{am})$$

$$P(\langle /S \rangle | \text{not})$$

$$= \frac{C(\langle S \rangle | I)}{C(\langle S \rangle)} \frac{C(I | I)}{C(I)} \frac{C(I | \text{am})}{I(I)} \frac{C(\text{am} | \text{not})}{C(\text{am})}$$

$$\frac{C(\text{not} | \langle /S \rangle)}{C(\text{not})}$$

$$= \frac{3}{3} \times \frac{1}{4} \times \frac{2}{4} \times \frac{1}{2} \times \frac{0}{1}$$

$$= 0$$

Q10 consider following training data.

$\langle S \rangle I \text{ am Jack } \langle IS \rangle$

$\langle S \rangle \text{ Jack I am } \langle IS \rangle$

$\langle S \rangle \text{ Jack I like } \langle IS \rangle$

$\langle S \rangle \text{ Jack I do like } \langle IS \rangle$

$\langle S \rangle \text{ do g like Jack } \langle IS \rangle$

Assume that we use a bigram language model based on above data.

10 what is most probable next word predicted by the model ?

1. $\langle S \rangle \text{ JACK } \dots$ 2. $\langle S \rangle \text{ Jack I do } \dots$

3. $\langle S \rangle \text{ JACK I am Jack } \dots$

4. $\langle S \rangle \text{ do g like } \dots$

15

$$\Rightarrow P(I | \langle S \rangle) = C(\langle S \rangle | I) / C(\langle S \rangle) = 1/5$$

$$P(\text{Jack} | \langle S \rangle) = C(\langle S \rangle | \text{Jack}) / C(\langle S \rangle) = 3/5$$

$$P(\text{do} | \langle S \rangle) = C(\langle S \rangle | \text{do}) / C(\langle S \rangle) = 1/5$$

$$P(\text{am} | I) = C(I | \text{am}) / C(I) = 2/5$$

$$P(\text{like} | I) = C(I | \text{like}) / C(I) = 2/5$$

$$P(\text{do} | I) = C(I | \text{do}) / C(I) = 1/5$$

$$P(\langle IS \rangle | \text{Jack}) = C(\text{Jack} | \langle IS \rangle) / C(\text{Jack}) = 2/5$$

$$P(\langle IS \rangle | \text{like}) = C(\text{like} | \langle IS \rangle) / C(\text{like}) = 2/3$$

$$P(\langle IS \rangle | \text{am}) = C(\text{am} | \langle IS \rangle) / C(\text{am}) = 1/2$$

$$P(I | \text{Jack}) = C(\text{Jack} | I) / C(\text{Jack}) = 3/5$$

$$P(\text{like} | \text{do}) = C(\text{do} | \text{like}) / C(\text{do}) = 1/2$$

$$P(I | \text{do}) = C(\text{do} | I) / C(\text{do}) = 1/2$$

$$P(\text{Jack} | \text{like}) = C(\text{like} | \text{Jack}) / C(\text{like}) = 1/3$$

$$P(\text{Jack} | \text{am}) = C(\text{am} | \text{Jack}) / C(\text{am}) = 1/2$$

Page : 28
Date :

1. < s > Jack I

2. < s > Jack I do like & I

3. < s > Jack I am Jack I

4. < s > do g like < /s >

Syntax Analysis

5. Syntax Analysis or parsing is the important phase in Natural Language Processing.

10. The purpose of this phase is to draw exact meaning or dictionary meaning from the text.

15. Syntax refers to the arrangement of words in a sentence such that they make grammatical sense.

20. In Natural language Processing Syntactic analysis is used to assess how the Natural language Aligns with the grammatical rules.

25. Computer Algorithms are used to apply grammatical rules to a group of words and derive meaning from them.

30. Syntactic analysis helps us understand the roles played by different words in a body of text.

Example: I am a good boy makes more sense
than am I bog a good.

Q1. Write a short note on Part of Speech Tagging?

=> The part of speech tagging is a process of assigning corresponding part of speech like noun, verb, adverb, etc. to each word in a sentence. It is a process of converting a sentence to forms - list of words, list of tuples (where each tuple is having a form (word, tag)). The tag in case of is a part of speech tag, and signifies whether the word is a noun, adjective, verb and so on.

- o The main challenge in POS tagging is to resolve the ambiguity in possible POS tags for a word. The word bear in the following sentence have same spelling but different meanings. She saw a bear · your efforts will bear fruit.
- o Hence it is impossible to have a generic mapping for Postags.
- o New type of contexts and new words keep coming up in dictionaries in various languages, and manual POS tagging is not Scalable in itself
- o That is why we rely on Machine based POS tagging.

Why Part of Speech Tagging

- o POS Tagging in itself may not be the solution to any particular problem of NLP
- o It is however something that is done as a prerequisite to simplify a lot of different problems.

Why is tagging hard.

- o Example

1. Book / VB that / DT flight / NN
 2. Does / VBZ that / DT flight / NN serve / VB dinner / NN
- Tagging is a type of disambiguation
Book can be NN or VB

- Can I read book on this flight?
- That can be a DT or complementizer
- My travel agent said that there would be a meal on this flight.
- - Words often have more than one word class: this
This is a nice day: PRP
This day is nice: DT
You can go this far: RB

VB - Verb - base form

DT - Determiner

NN - Noun, singular or mass

VBZ - Verb, 3rd person singular present.

PRP - Personal Pronoun

RB - Adverb.

VBP - Verb, non 3rd person singular present.

Applications of PoS Tagging in various NLP tasks :

1. Text to speech conversion:

"They refuse to permit us to obtain the refuse permit"

o The word refuse is being used twice in this sentence and has two different meanings here.

o 'refUSE' is a verb meaning 'deny' while 'REFuse' is a noun meaning 'trash'

o PoS tags generated for this sentence by the NLTK package.

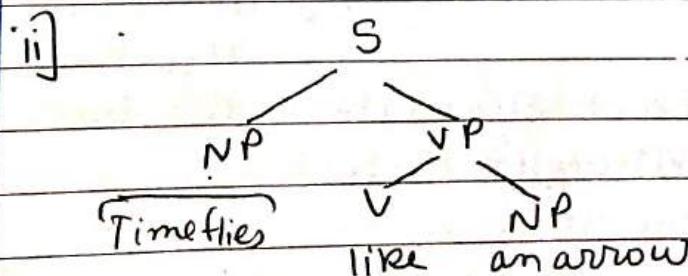
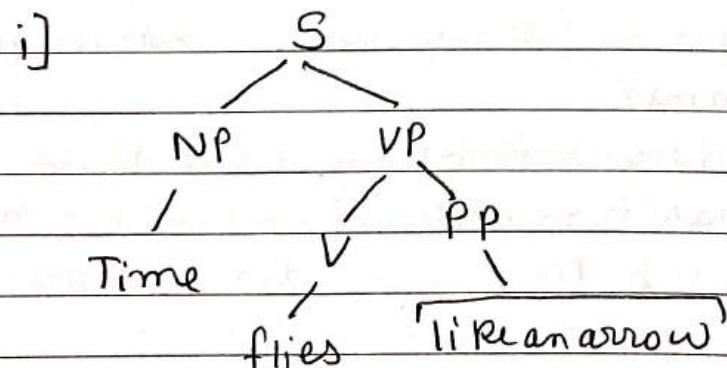
>>> text = word_tokenize ("They refuse to permit us to obtain the refuse permit") >>>
 nlpR.pos_tag (text) [('They', 'PRP'), ('refuse', 'VBP'),
 5 ('to', 'TO'), ('Permit', 'VB'), ('us', 'PRP'), ('to', 'TO'),
 ('obtain', 'VB'), ('the', 'DT'), ('refuse', 'NN'),
 ('permit', 'NN')]

- o POS tags for both refuse and REFuse are different.
- o using these two different pos tags for our text to speech converter can come up with a different sound.

2. Word Sense Disambiguation (WSD)

- o WSD is identifying which sense of a word (that is, which meaning) is used in a sentence, when the word has multiple meanings.
 consider a Sentence.

"Time flies like an arrow"



* Types of POS Taggers:

POS Tagging algorithms fall into two distinctive groups:

- o Rule based PostTaggers.

- o Stochastic PostTaggers.

i] Rule - Based POS Tagger

- o They use dictionary or lexicon for getting possible tags for tagging each word.

- o If the word has more than one possible tags, then rule based tagger uses hand-written rules to identify the correct tag.

- o Disambiguation can also be performed in rule based tagging by analyzing the linguistic features of a word along with its preceding as well as following words.

- o For example suppose if the preceding word of a word is article then the word must be a noun.

- o All such kind of information in rule based POS tagging is coded in the form of rules.

- o These rules may be either:

- Context - pattern rules

- or as Regular Expression compiled into Finite State Automata, intersected with lexically ^{seected}.

ambiguous sentence representation.

- o Rule-Based POS Tagging has a two stage architecture:

5 o first stage: It uses a dictionary to assign each word a list of potential parts of speech

10 o Second Stage: It uses large lists of handwritten disambiguation rules to sort down the list to a single part of speech for each word.

Properties of Rule-Based POS Tagging

15 o These taggers are knowledge driven taggers

20 o The rules in Rule based POS tagging are built manually

o The information is coded in the form of rules

25 o We have some limited number of rules approximately around 1,000

30 o Smoothing and language modeling is defined explicitly in Rule based taggers.

ii] Stochastic POS Tagger :

- o The model that includes frequency or probability can be called stochastic.
- o Any number of different approaches to the problem of part of speech tagging can be referred as stochastic tagger.

Approach for POS Tagging (Stochastic) :

i] Word Frequency Approach :

- Stochastic taggers disambiguate the words based on the probability that a word occurs with a particular tag.

The tag encountered most frequently in the training set is the one assigned to the ambiguous instance of the word.

ii] Tag Sequence Probability

- Here, the tagger calculates the probability of a given sequence of tags occurring

- The best tag for a given word is determined by the probability that it occurs with the n previous tag.

- It is also called n -gram approach.

Properties of stochastic Pos Tagging

It is called so because the best tag for given word is determined by the probability at which it occurs with one n previous tags.

Properties of stochastic Pos Tagging.

- o This Pos Tagging is based on the probability of tag occurring.
- o It requires a training corpus.
- o There would be no probability for the word that do not exist in the corpus.
- o It uses different testing corpus.
- o It uses the simplest Pos tagging because it chooses most frequent tags associated with a word in training corpus.

3. Transformation Based Tagging :

- o Transformation based tagging is also called Brill tagging.
- o It is a instance of transformation based Learning, which is a rule based algorithm for automatic tagging of POS to the given text.
- o Like Rule based, it is also based on the rules that specify what tags need to be assigned to what words.
- o Like Stochastic, it is machine Learning technique in which rules are automatically induced from data.
- o Hence transformation Based tagger draws inspiration from both rule based & Stochastic taggers.

Working of Transformation Based Tagging:

- o To understand TBT, we need to know the working of TBL (Transformation Based Learning)
 - Start with the solution: It usually starts with some solution to one problem and works in cycles.

- Most beneficial transformation chosen:
In each cycle, TBL will choose the most beneficial transformation

5 - Apply to the problem:

The transformation chosen in the last step will be applied to the problem.

- The algorithm will stop when the selected transformation in step two will not add either more value or there are no more transformations to be selected.

15 • Advantages of Transformation-Based Learning

- We learn small set of simple rules and these rules are enough for tagging.

- Development as well as debugging is very easy in TBL because the learned rules are easy to understand.

- Complexity in tagging is reduced because in TBL there is interlacing of machine learned & human generated rules.

- Transformation based tagger is much faster than Markov Model tagger.

Page : 40
Date :

Disadvantages of Transformation Based learning

- TBL does not provide tag probabilities
- ⁵ In TBL the training time is very long especially on large corpora.

Q Write a short note on multiple tags, multiple words and unknown words.
Multiple tags & Multiple words.

=> Two issues that arise in tagging are tag indeterminacy and multi-part words. Tag indeterminacy arises when a word is ambiguous between multiple tags and it is impossible or very difficult to disambiguate. In this case, some taggers allow the use of multiple tags. This is the case in the Penn Treebank and in the British National Corpus. Common tag indeterminacies include adjective versus preterite versus past participle (JJ / VBD / VBN), and adjective versus noun as prenominal modifier (JS / NN).

The second issue concerns multi-part words.

The C5 and C7 tagsets, for example, allow prepositions like 'in terms of' to be treated as a single word by adding numbers to each tag:

in/1/131 terms/1/132 of/1/133

Finally, some tagged corpora split certain words; for example the Penn Treebank and the British National Corpus splits contractions and the 's-genitive from their stems:

would/MD n't/RB

children/NNS · 's/POS

Unknown Words :

This concept is based on the idea that if we consider a language for example English so everyday there are many new words added or created in this language which don't have any part of speech.

So this concept is based on how to assign a part of speech to that word.

There are 3 strategies or ways.

1. Ambiguity Declaration: Declare that word as an ambiguous word and it has every part of speech (noun, verb, adjective etc) and with the help of trigram get rid of this ambiguity. Trigram means two words before that word or two words after that word or both, and calculate try to understand the context with the help of this to get rid of the ambiguity.

example if a word a is before a word so the next word would be noun. or if the previous word is to the next word would be verb.

2. Spelling check:

This concept is totally based on checking the spelling of a word which has several rules. for example if the word has ed after it then the word is verbs past participle eg cooked. if a word starts with capital letter its a noun. In this way the POS tag is assigned to that word.

Page : 43
Date :

3. Transformation Based Learning :

- Assigning a POS Tag to a word by the combination of various techniques eg length of word, starting, ending, context of the word.

5

10

Context Free Grammar

A Context Free Grammar is a list of rules that define the set of all well formed sentences in a language. Each rule has a left hand side, which identifies a syntactic category, and a right hand side, which defines its alternative component parts, reading from left to right.

There are three main concepts: constituency, grammatical relations, and subcategorization and dependencies.

15 Constituency: The fundamental idea of constituency is that groups of words may behave as a single unit or phrase called a constituent. For example we will see that a group of words ~~can~~ called a noun phrase often acts as a unit; noun phrases include single words like She or Ojas and phrases like the house, Russian girl and a well-weathered three-story structure.

25 Grammatical Relations: Grammatical Relations are formalization of ideas from traditional grammar about SUBJECT and OBJECT. In the sentence: She ate a mammoth breakfast. The noun phrase She is the SUBJECT and a mammoth breakfast is the OBJECT.

Page : 95
Date :

Subcategorization & Dependency:

Subcategorization and Dependency relations refer to certain kinds of relations between words and phrases.

For example the verb want can be followed by an infinitive, as in I want to fly to Detroit, or a noun phrase, as in I want a flight to Detroit.

But the verb find cannot be followed by an infinitive (*I found to fly to Dallas). These are called facts about the subcategory of the Verb. All of these kinds of syntactic knowledge can be modelled by various kinds of grammars that are based on context free grammars.

Context free grammars are thus the backbone of many models of the syntax of natural language.

Q) Construct a parse tree for the following sentence using CGF rules:

"The man read this book"

rules: $S \rightarrow NP VP$

$S \rightarrow AUX NP VP$

$S \rightarrow VP NP \rightarrow Det N OM$

$N OM \rightarrow N OM$ $VP \rightarrow V erb N P$

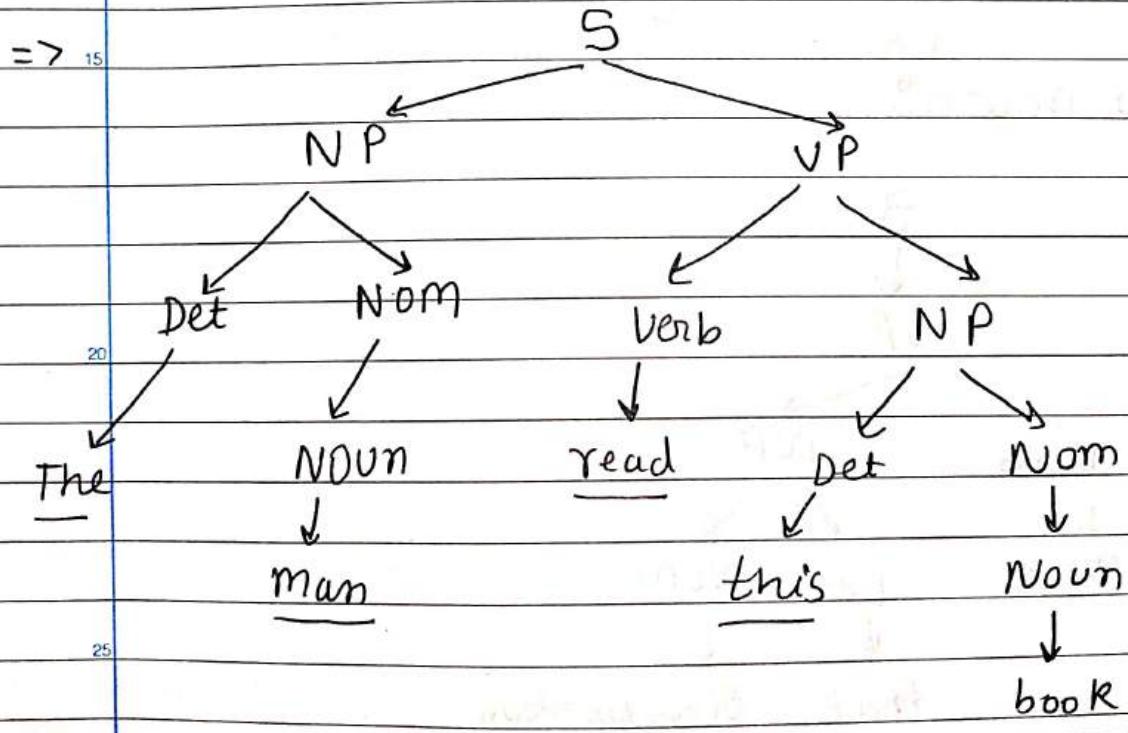
$N OM \rightarrow N OM N OM$

$Det \rightarrow t hat / t his / a / t he$

$N OM \rightarrow b ook / f light / m eal / m an$

$V erb \rightarrow b ook / i nclude / r ead$

$A ux \rightarrow d oes$



Q

Generate the input sentence by the help of given grammar [Top down and Bottom up]

5 Grammer:

$S \rightarrow VP$

$VP \rightarrow \text{verb } NP$

$NP \rightarrow \text{Det } Noun$

$ND \rightarrow \text{Det } Noun$

10 $\text{Det} \rightarrow \text{that}$

$\text{Noun} \rightarrow \text{Singular Noun}$

$\text{Verb} \rightarrow \text{Book}$

$\text{Singular noun} \rightarrow \text{flight}$

input $\rightarrow \text{Book that flight.}$

15

Top Down:

20 S

↓

VP

25 Verb NP

↓

Book

→

Det

Noun

↓

↓

30 that Singular Noun

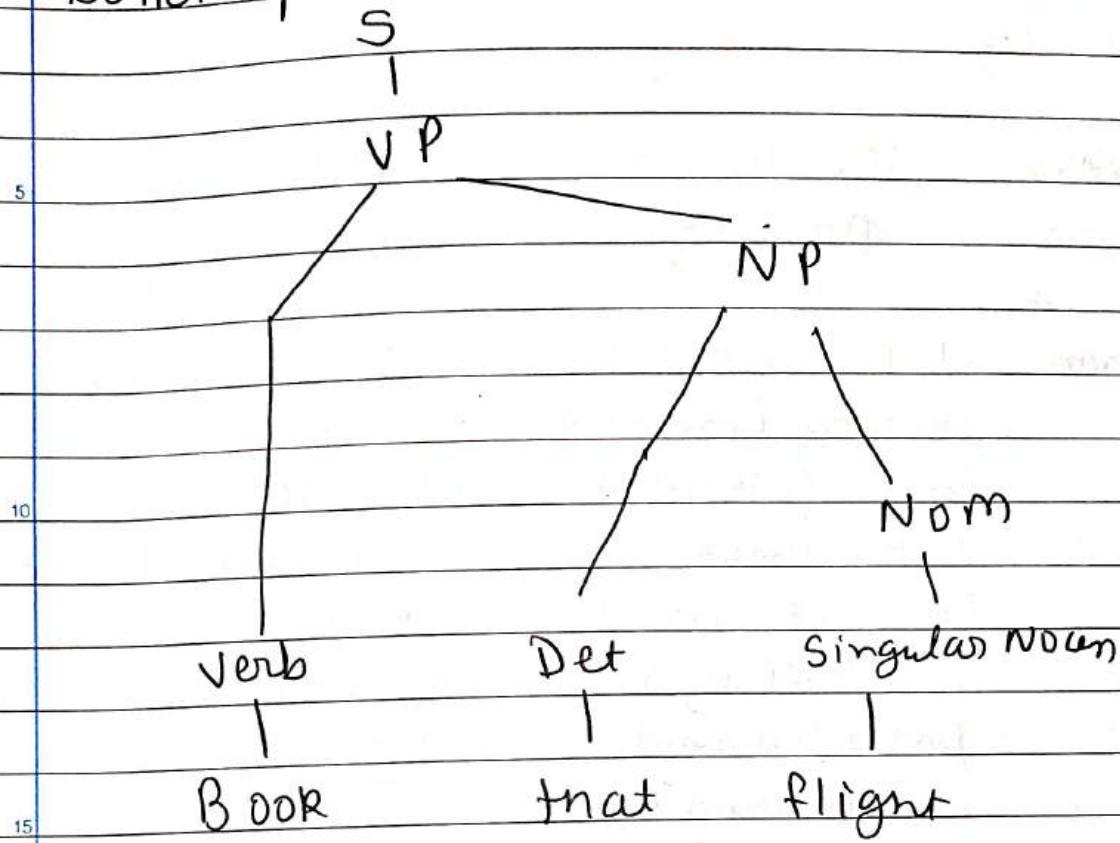
↓
flight.

25

30

Page : 48
Date :

Bottom up



Q

Differentiate Between top down and Bottom up parsing.

Parameters	Top Down parsing	Bottom Up parsing
------------	---------------------	----------------------

Definition	It performs the parsing from one starting symbol to the input string. It starts from the root level of the parse tree and works down by using the rules of formal grammar.	A parsing strategy that first looks at the level of the parse tree and works up the parse tree by using the rules of a formal grammar.
Strength	Moderate strength	Powerful than Top down parsing
Main Decision	To Select what production rule to use in order to construct a string	To select when to use a production rule to reduce the string to get the starting symbol
Method of construction	Left most Derivation	Right most Derivation
Example	Recursive Descent Parser	SR Parser.

Semantic Analysis

What is Semantic Analysis?

- => As a human being making sense of text is very simple we recognize every word and the context in which it is used.
- For example if you read this sentence "Your teaching is very good, I scored full marks in my Viva".
- So here you understand that student is very happy with the teaching skills of the teacher because he scored full marks in his Viva.
- However machines first need to be trained to make sense of the human language and understand the context in which words are used or else they might misinterpret the word good as negative.
- With the power of Machine Learning algorithms and Natural Language Processing Semantic analysis systems can understand the context of natural language, detects emotions and sarcasm and extract valuable information from unstructured data, achieving human-level accuracy.

Page : 51
Date :

- In a very simple terms we can say Semantic Analysis is the process of finding out meaning from the text.
- The word Semantics means the study of meaning . Semantic Analysis allows computers to understand and interpret sentences , paragraphs or whole documents , by analyzing their grammatical structure , and by identifying relationships between individual words in a particular context .
- In other or technical terms we can say that semantic analysis is the process whereby meaning representations are composed and assigned to linguistic inputs .

Applications of semantic analysis

1. Information Extraction

Information Extraction is a precise process of extracting information from unstructured textual sources to enable finding entities as well as classifying and storing them in the database.

2. Text Summarization

Text Summarization can automatically shorten longer texts and extract summaries of sections of text without losing the message.

3. Information Retrieval System

Information Retrieval System is the process of tracing and recovering specific information from stored data.

4. Machine Translation

Machine Translation is the process of translating one source language text into another language and it is one of the most important applications.

Page : 53
Date :

5. Expert system

Expert system is a system which gives you expert advice.

For example : Grammarly.

Wordnet

- WordNet is a huge lexical database of English.
- Nouns, verbs, adjectives and adverbs are grouped into sets of cognitive Synonyms & each expressing a distinct concept.
- These Synonyms are called as Synsets and they are interlinked by means of conceptual-
 - Semantic and lexical relations.
- The resulting network of meaningfully related words and concepts can be navigated with the browsers.
- WordNet is freely available on Internet, it is public and anyone can download it and make use of it.
- The structure of WordNet make it a useful tool for computational linguistics and ~~narration~~
NLP

Structure of WordNet

- The main relation among words in WordNet is Synonymy, as between the words shut and close or car and automobile.
- Synonyms - words that denote the same concept and are interchangeable in many contexts - are grouped into unordered sets (synsets).
- There are ~~various~~ many unique forms of different categories in WordNet.

Category	Unique Forms
Noun	117,097
Verb	11,488
Adjective	22,141
Adverb	4,609

- Each of WordNet's synsets is linked to other synsets by means of a small number of "conceptual relations". Additionally, a synset contains a brief definition ("gloss") and, in most cases, one or more short sentences illustrating the use of the synset members.
- Words forms with several distinct meanings are represented in as many distinct synsets.
- Thus each form-meaning pair in WordNet is unique.

Applications of WordNet

- WordNet has lot of Applications in IR and NLP (IR - Information Retrieval)
Some of applications are mentioned below

1. Automatic Query Expansion:

WordNet Semantic Relations can be used to expand queries so that the search for a document is not confined to the pattern-matching of query terms, but also covers Synonyms.

2. Document Summarization:

WordNet has found useful application in text summarization. Few approaches utilize information from WordNet to compute lexical chains.

3. Concept Identification in Natural Language:

WordNet can be used to identify concepts pertaining to a term, to suit them to the full & semantic richness and complexity of a given information need.

4. Word Sense Disambiguation

WordNet.com combines features of a number of the other resources commonly used in disambiguation work. It offers sense definitions of words, identifies synsets

of synonyms, defines a number of semantic relations and is freely available.

This makes it the best known and most

Page : 57
Date :

~~Used~~ used resource for word sense disambiguation.

Word Sense Disambiguation

=> Many of times a single word might have multiple meanings.

- Word Sense Disambiguation in Natural Language Processing is defined as the ability to determine which meaning of word is activated by the use of word in a particular context.
- For example there is a sentence

She killed him with a baseball bat.

In this sentence one meaning of bat the word "bat" is it is a mammal and the other meaning of the word "bat" is it is a wooden object.

- Now if we consider this statement the word bat has multiple meaning.

- So the process of identifying the correct meaning of such a word according to the statement in which the word is, is called Word sense disambiguation.

- Now word sense disambiguation has different approaches
 - 1. Knowledge based approach or Dictionary approach
 - 2. Supervised and Unsupervised approach
 - 3. Hybrid Approach.

Knowledge based Approach

- o Requirements:

- 1. Raw corpora
- 2. Machine understandable / readable dictionary eg: Indoword.

How knowledge based Approach Works

There is a creation of two bag of sense bag and a context bag

Consider this sentence and the meaning of the word "bank" in the dictionary wordnet below

Sentence: The bank can guarantee deposits will eventually cover future tuition costs because it invests in adjustable rate mortgage.

The meaning / senses of the word bank in the wordnet dictionary are

- 5 Bank sense⁰: A financial institution that accepts deposits and channels the money ^{into} lending activities.
 Example: "That bank holds the mortgage on my home"

10 Bank² sense: Sloping land (especially beside a body of water)

Example: They pulled the canoe up on the bank.

- In Knowledge based approach there are two bags sense bag and context bag
- Sense bag consists of all the senses of the word bank (because bank is the word with multiple meanings).
- And context bag consists remaining content other than the word "bank" from the sentence.
- Now context bag is compared with each of the sense of the word bank.

- The sense from which maximum words are matched with the sentence is chosen as the meaning of the word "bank".
- In this case sense 2 doesn't have any word matching with any word from context bag.
- But in the case of sense 1 there are two words (Deposit and Mortgage) which matches with words from context bag.
- So the Sense 1 is the meaning of the word bank.

Page : 62
Date :

LesR Algorithm

- LesR Algorithm is a very classical algorithm for word sense disambiguation.
~~LesRA~~
- LesR Algorithm is a knowledge based approach.
- LesR Algorithm is an classical algorithm that is used for word Sense Disambiguation.

Page : 63
Date :

Homonymy

- Homonyms are those words that have same spelling and pronunciation but when it comes to their meaning, their meaning is different.

for example:

i) bat (wooden object)

vs bat (flying weird mammal)

ii) bank (riverbank or riverside)

vs bank (financial institution)

Polysemes

- Polysemous words are words with different but related meanings
- A word will become polysemous if it expresses different meanings
- Where the difference between meanings can be subtle or obvious
- If we examine the origins of words it can help to decide whether a word is polysemous or homonymous. This might need to be done because it is sometimes difficult to make out whether a word is polysemous or not because the connection between words can be unclear.

example :

1. He drank a glass of cho

2. The angry artist sued the newspaper

3. She read ^{the} newspaper .

Synonymy

- ~~to~~ Synonyms are the words with similar meanings.
- Some synonyms don't have exactly the same meaning - there might be a tiny difference.

Example :

1. Beautiful - Gorgeous
2. Mistake - Error
3. Rich - Wealthy.

Antonymy

- Words with opposite or contrasting meanings ~~are~~ share the relationship called as Antonymy
- There are three types of Antonyms ~~which~~ which are complementary, gradable and relational antonyms.
- Complementary antonyms are pairs of words which have opposite meanings that do not lie on a continuous spectrum
Example : interior - exterior
right - wrong
Exhale - inhale etc.

- Graddable antonym are pairs of words with opposite meanings that lie on a continuous spectrum

5 Example if we take age as a continuous spectrum, young and old are two ends of the spectrum.

-₁₀ Relational antonyms are the pairs of words that refer to a relationship from opposite points of view

- Example : Patient : doctor
15 brother : sister etc.

20

25

30

Hypernymy and Hyponymy

→ Hyponymy

It is a sense which is a subclass of another sense

i) Car is a hyponym of vehicle

ii) dog is a hyponym of animal

iii) strawberry is a hyponym of fruit

→ Hypernym is a sense which is a superclass

i) animal is a hypernym of dog

ii) fruit is a hypernym of strawberry

iii) vehicle is a hypernym of car

→ Hyponymy is a transitive relation,
if a is hyponym of b, b is hyponym of c
then a is hyponym of c

- for example if violet is hyponym of purple and purple is hyponym of color
then violet is hyponym of color.

- A word can be both a hypernym and hyponym : for example purple is a hyponym of color but itself is a hypernym of its shades.

Page :
Date :

68

Meronymy

- A meronym is a word which represents a constituent part of something
- if x is a part of y then x is meronym of y
- for example coconut is a meronym of coconut tree.

Pragmatics

Pragmatics & Discourse

5 Pragmatics:

Concerns how sentences are used in different situations and how it affects the interpretation of the sentence.

10

Discourse:

Concerns how the immediately preceding sentences affect the interpretation of the next sentence

15

Basic difference:

20 Pragmatics analyses individual utterances (organized set of words) in context. Discourse focusses on an organized set of utterances.

25

30

Pragmatics:

- Extension of semantics and proposition logic.
- Studies meaning of utterance and ~~code~~ defines rules to govern their interpretation.
- Difference between semantics and Pragmatics is semantics: Literal meaning
Pragmatics: Intended meaning
(needs context
Information)
- Pragmatics uses context of Utterance
 - When, why, by who, to whom...
Something is said.
- Also it deals with intentions
 - Criticize, inform, promise, request, warn...
- There are three major applications of pragmatics
 1. Question Answer System
 2. Summarization
 3. Sentiment Analysis.

Discourse

- It is a study of meaning too, but it focusses on large scale units (Articles, conversations...) and overall interpretation in a specific context.
- Discourse means the group of related Sentences
 - for example
 - Prime Minister's Speech: Namaste! Mere deshwasiyo.
here Semantic Analysis shall be used
 - 'Ladies' word beside the seat in public transport
here pragmatic analysis shall be used.
 - Namaste! Mere deshwasiyo. Ghar se bahar mat niklo.
here discourse analysis shall be used
 - Discourse analysis involves the study of the relationship between language and contextual background.

Now here Contextual background include

- Situational context -

knowledge about physical situations existing in the surroundings at the time of utterance.

- Background knowledge :

includes cultural knowledge and interpersonal knowledge.

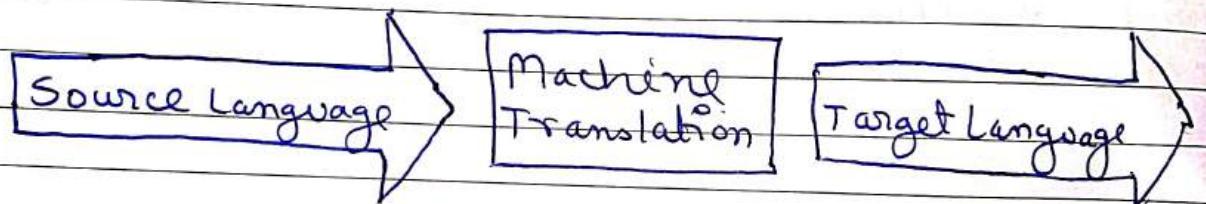
- Co-textual context -

knowledge of what has been said earlier.

Applications

Machine Translation

Machine Translation (MT) is the automated process of translating one natural language to another. Machine translation, an integral part of Natural Language Processing where translation is done from source language to target language preserving the meaning of the sentence.



-20 There are few challenging aspects of MT:

1) The wide variety of languages, alphabets and grammar;

2) The task to translate a sequence to a sequence is harder for a system than working with numbers only

3) There is no one correct answer

- Different types of Machine Translation

1 Statistical Machine Translation (SMT)

- SMT functions by referring to statistical models that are based on the analysis of large volumes of bilingual text
- It ~~also~~ work towards or aims to determine the correspondence between a word from the target language and a word from the source language.
- Google translate is an example!
- SMT is good ~~as~~ ^{for} basic translation. But its disadvantage is that it does not factor in context, which means translations can often be erroneous. We can also say that it doesn't expect high quality translations

2. Neural Machine Translation (NMT)

- It is a new approach that makes machines learn to translate through one large neural network (multiple processing devices modeled on the brain)
- The approach has become increasingly popular among MT researchers and developers, as trained NMT Systems have begun to show better translation performance in many language pairs compared to the

phrase based statistical approach.

3. Rule based Machine Translation (RBMT)

- It translates on the basis of grammatical rules
- To generate the translated sentence RBMT conducts a grammatical analysis of the source language and the target language.
- It requires proof reading, and its heavy dependence on lexicons means that efficiency is achieved after a long period of time

4. Hybrid Machine Translation (HMT)

- It is a mix of RBMT and SMT. It takes up, the translation memory making it far more effective in terms of quality
- HMT ~~is~~ also has its disadvantages:
 1. It needs ~~is~~ heavy editing
 2. There is a requirement of Human Translators.

Information Retrieval (IR) or

Write a Short Note on Information Retrieval (IR)

so): Information Retrieval is one of the most challenging problems of Natural Language

- IR is defined as a software program that deals with the organization, storage, retrieval and evaluation of information from document repositories particularly textual information.

- IR system assists users in finding the information and it informs the existence and location of documents that might consist of the required information.

- The documents that satisfy the user's requirement are called relevant documents.

Example : Google, Yahoo, Altavista etc.

- Traditionally, the IR System techniques are based on keyword.

- They use lists of keywords to describe the content of information but they do not say reveal semantic relationships between keywords nor consider the meaning of words and phrases.

=> Basic IR system involves following Stages:

1. Indexing the collection of documents.
2. Transforming the query in the same way as the document content is represented
3. Comparing the description of each document with that of the query.
4. Listing the results in order of relevancy

- In general all IR Systems consists of mainly two processes as

- a) Indexing : Indexing is the process of selecting terms to represent a text which involves tokenization of String, removing frequent words and stemming.
- b) Matching : It is the process of computing a measure of similarity between two text representations. Relevance of a document is computed based on parameters like term frequency and inverse document frequency.

Question Answers System

Question Answering (Q A) System is a task of automatically answering to the questions asked in natural language using either a pre structured database or a collection of natural language documents

- 10 It presents only the requested information instead of searching full documents like Search engine.
- The basic idea behind the Q A System is that the users just have to ask the question and the system will retrieve the most appropriate and correct answer for the question

20 Example

Q. "What is the birth place of Shree Krishna"?

A. Mathura.

- 25 Question answering System helps users to find the precise answers to the question articulated in natural language.
- 30 Question answering system provides explicit, concise and accurate answer to user questions rather than providing a set of relevant documents or web pages as

Answers as most of the information retrieval system does.

- o Question Answering System basically consists of three parts as - Question processing - answer retrieval - answer generation.
- o QAS has become part of daily life of users.
- o Over a period of time many personal assistance software like Siri, Cortana, Google Now, Alexa etc. are developed which provide precise and accurate answer to user's questions.
- o Datasets for Q A Systems are
 1. Stanford Question Answer Dataset
 2. Wiki QA dataset
 3. TREC-QA
 4. News - QA

=> Question Answering system challenges

- Lexical Gap: In a natural language, the same meaning can be expressed in different ways. Because a question can usually only be answered if every referred concept is identified, bridging this gap significantly increases the proportion of questions that can be answered by a system.

- Ambiguity : It is the phenomenon of the same phrase having different meanings. This can be structural and syntactic (like 'flying planes') or lexical and semantic ('like bank'). The same string accidentally refers to different concepts (as in money bank vs. river bank) and polysemy, where the same string refers to different but related concepts (as in bank as a company vs bank as a building)
- Multilingualism : Knowledge on the Web is expressed in various languages. While RDF resources can be described in multiple languages at once using language tags, there is not a single language that is always used in web documents.
- Additionally : Users have different native languages. A Q&A system is expected to recognize a language and get the results on the go!

Text Categorization System

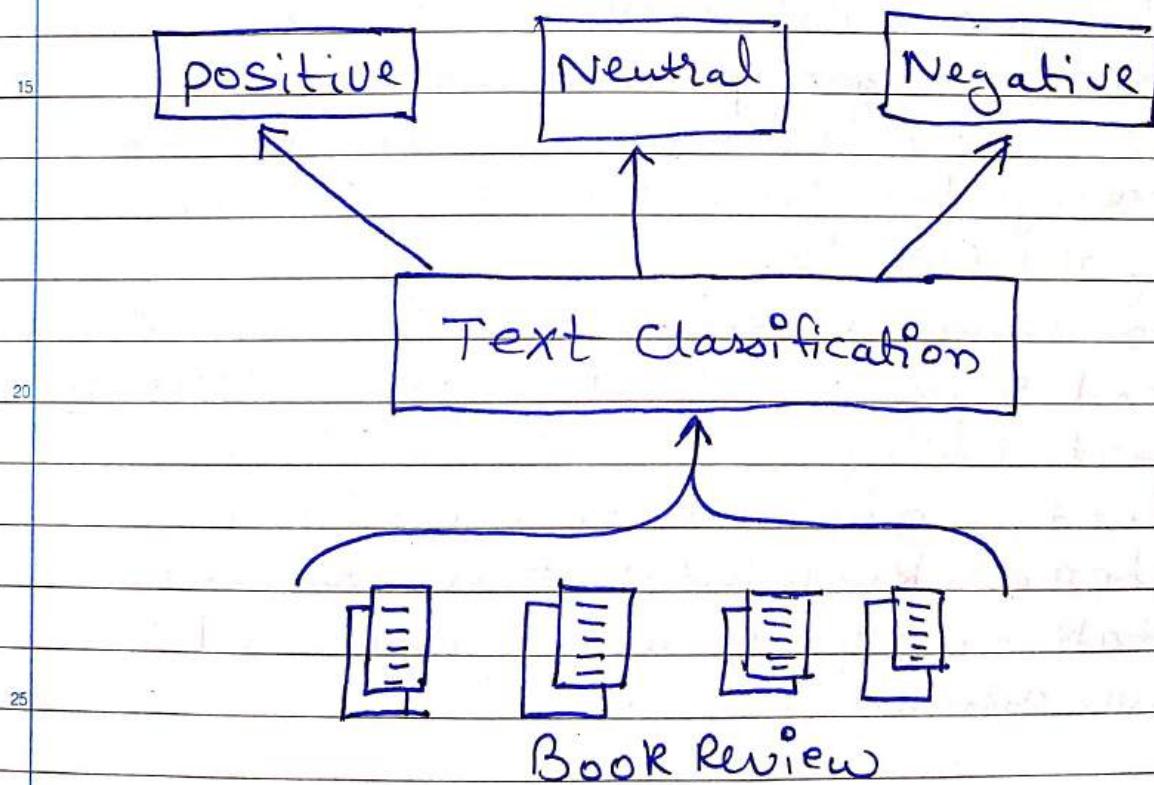
Need of text categorization System

- o Rapid development of Information technology has led to massive growth in text data found on internet.
- o Data mining field worked mostly on English text document.
- o Nowadays, millions of documents are present in Indian ~~and~~ regional languages like Telugu, Tamil, Hindi, Punjabi, Bengali, Urdu, Marathi.
- o To classify such documents manually is an expensive and time consuming task.
- o Automatic classification can help in better management and retrieval of these text documents.
- o Also their accuracy and time efficiency is much better than manual text classification.

- Text categorization (also known as text classification or topic spotting) is the task of automatically sorting a set of documents into categories (clusters)

o USES OF TEXT CATEGORIZATION

- Filtering of content
- Spam filtering Identification of document content
- Survey coding.



- Applications of text categorization
- Email message filtering
- News & event tracked and filtered by topics
- Web pages organized into categories hierarchy
- Text classification can be achieved through three main approaches
 - o Rule based approach: These app use of handcrafted linguistic rules to classify text is been made here.
 - A way to group text is to create list of words related to a certain column and then judge the text based on occurrences of these words.
 - for example ~~about~~ 'fur', 'feathers', 'claws' and 'scales' could help zoologist identify texts talking about animals online.
 - But this approach require a lot of domain knowledge to be extensive, take a lot of time to compile and are difficult to scale.
- o Machine learning approach: Machine learning is used to train models on large scale of text data to predict categories of new text. For the training of models, we need to transform text data into numerical data - this is feature extraction.

Important feature extraction techniques include bag of words and n-grams. There are many useful ML algos which can be used for text classification.

5 Naive Bayse classifiers, SVM are the most popular or famous once.

o Hybrid Approach

10 They are a combination of ML Approach and Rule-based approach. They make use of these two to model a classifier that can be fine-tuned in certain scenarios.

15 ⇒ Language detection, NLP, Topic detection, Semantic analysis are some of the most common examples and use cases for automatic text classification.

Text Summarization

- I don't want a full report, just give me a summary of the results. I have often found myself in this situation - both in college as well as my professional life. We prepare a comprehensive report and the teacher/supervisor only has time to read the summary.
- Summarization means to reduce the size of the document without changing its meaning.
- Automatic text summarization is the task of producing a concise and fluent summary while preserving key information content and overall meaning.
- A good summary should cover the most vital information of the original document or a cluster of documents, while being coherent, non-redundant and grammatically readable.

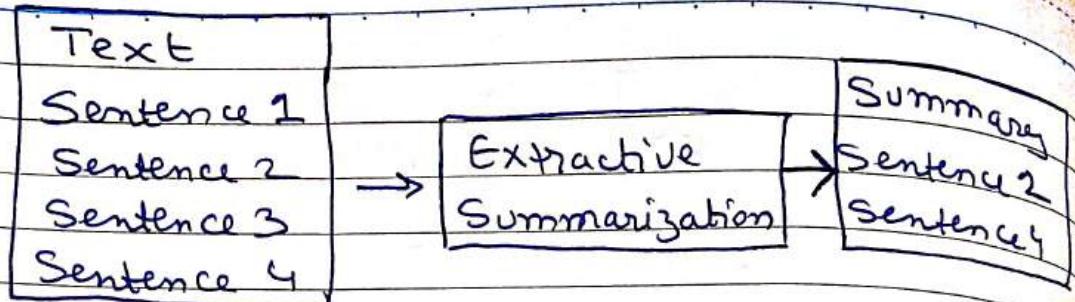
Text Summarization

5
Abstractive
Summarization

Extractive
Summarization

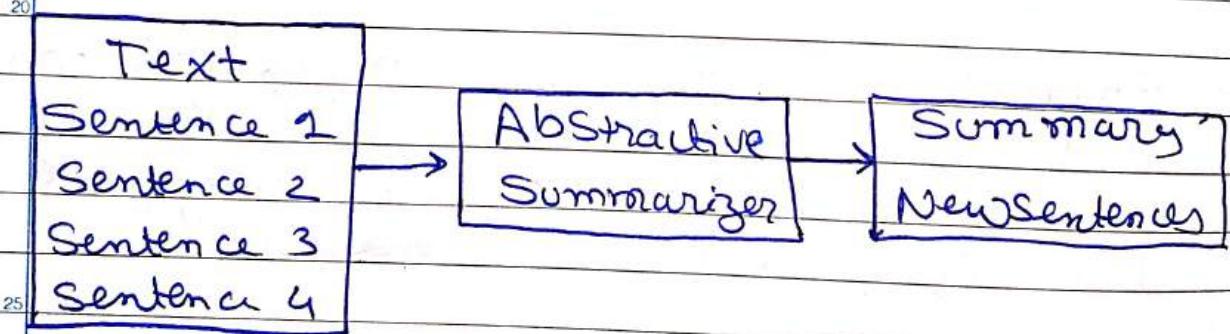
- 10 Extractive based Summarization

- o The extractive text summarization technique involves pulling keyphrases from the source document and combining them to make a summary.
- o The extraction is made according to the defined metric without making any changes to the texts.
- o Source text: Joseph and Marya rode on a donkey to attend the annual event in Jerusalem. In the city Marya gave birth to a child named ~~Jesus~~ Zeus
- o ~~Extr~~ Extractive summary: Joseph and Marya attend event Jerusalem. Marya birth Zeus.

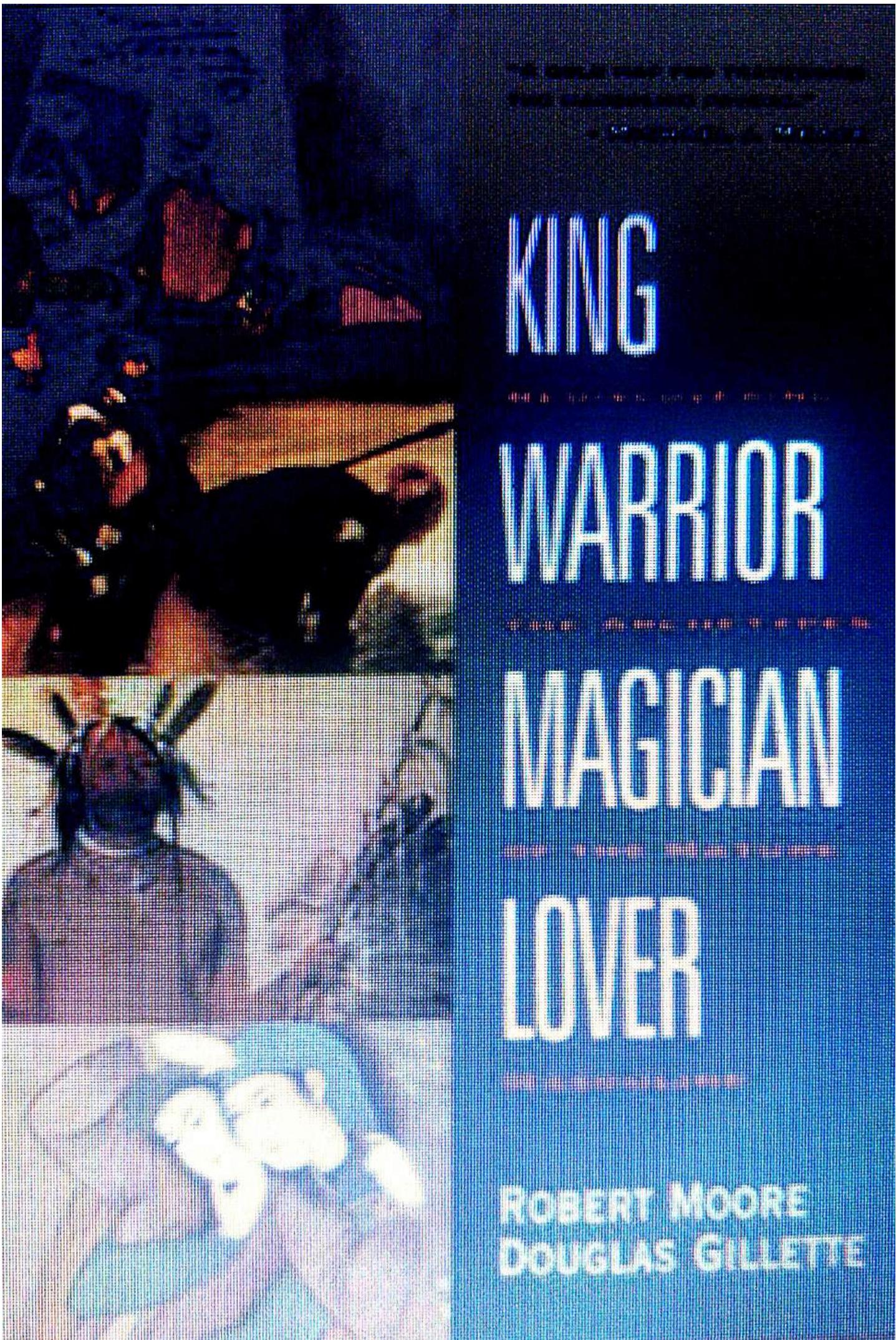


- Abstraction based Summarization

- o The abstraction technique entails paraphrasing and shortening parts of the source document.
- o The abstractive text summarization algorithms create new phrases and sentences that relay the most useful information from the original text. just like humans do. Therefore abstraction performs better than extraction.



- o Applications of text summarization
 - can be used as a preliminary stage for information retrieval tasks
 - simplifies text categorization
 - widely used due to information overload problem where information searched is very large and where there is a need of meaningful summary and saves time.



Tiny Changes, Remarkable Results



HABITS
FORMED
ARE
REFINED
BY
REPEATED
PRACTICE

An Easy & Proven Way to
Build Good Habits & Break Bad Ones

James Clear

INTERNATIONAL BESTSELLER

THE WAY OF THE SUPERIOR MAN

20th

Anniversary
Edition

*A Spiritual Guide To Mastering the Challenges
of Women, Work, and Sexual Desire*

DAVID DEIDA

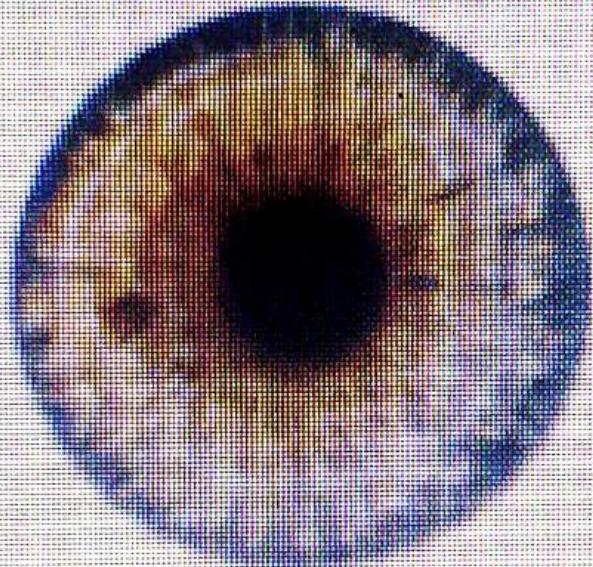
**MEN ARE
FROM MARS,
WOMEN ARE
from VENUS**

• JOHN GRAY •



FROM THE AUTHOR OF *SAPIENS*

Yuval Noah Harari



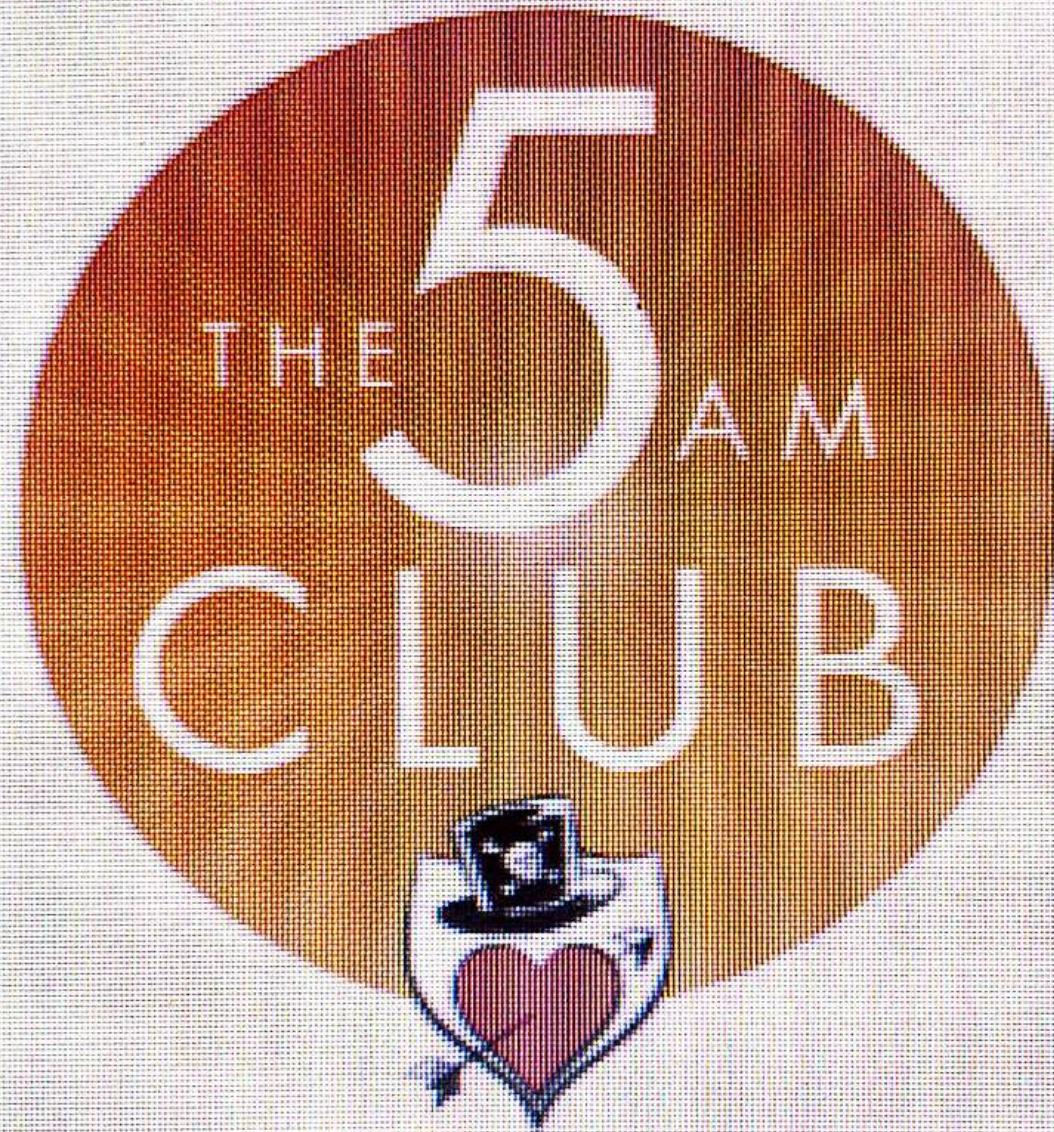
21 Lessons
for the
21st Century

Perfect Computer Engineer [YouTube]

THE #1 NEW YORK TIMES & #1 #1 NEW YORK TIMES BESTSELLING AUTHOR OF THE MONK WHO SOLD HIS FERRARI

ROBIN SHARMA

15 MILLION BOOKS SOLD WORLDWIDE



OWN YOUR MORNING
ELEVATE YOUR LIFE

The 48 LAWS OF
POWER

P O W E R

RICHARD
GREEN

Abridged Edition with
New Edition

JORDAN B. PETERSON

12 RULES
FOR LIFE

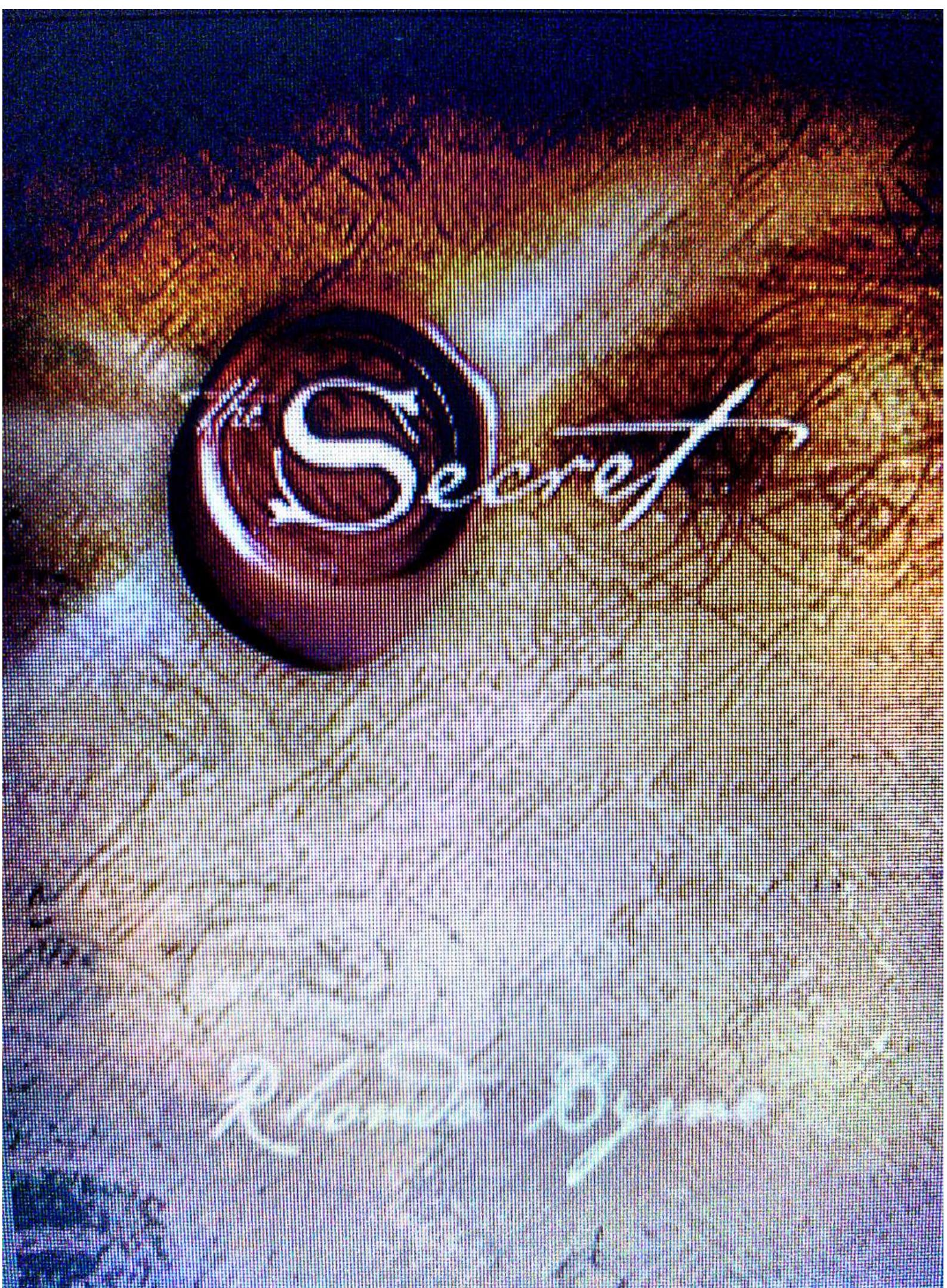
AN ANTIDOTE TO CHAOS

Foreword by Shilpa Shetty Kundra

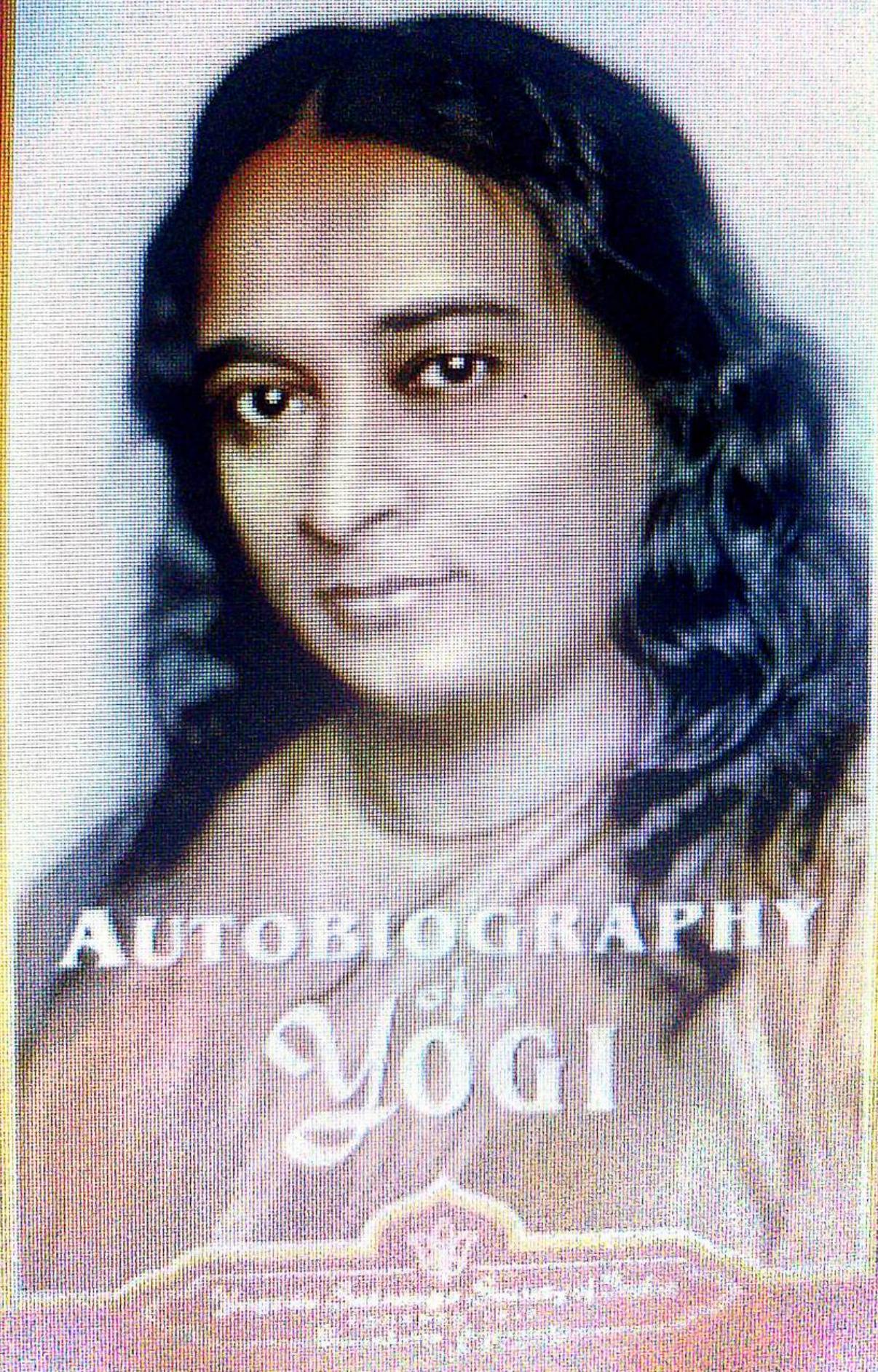
THE
Magic
WEIGHT-LOSS
PILL

62 lifestyle changes





Paramahansa Yogananda



THE
LAWNS OF POWER

THE
LAWS
OF
HUMAN
NATURE

THOMAS GREENE

Everyday ayurveda



Daily Habits That Can Change
Your Life In A Day

SHASHWATI CHATTACHARJEE

NEW YORK TIMES BESTSELLER

"*Sapiens* tackles the biggest questions of history and of the modern world, and it is written in unforgomable vivid language."

—JARED DIAMOND, Pulitzer Prize-winning
author of *Guns, Germs, and Steel*

Yuval Noah Harari



Sapiens

A Brief
History of
Humankind

THE NEW YORK TIMES BESTSELLER

Yuval Noah Harari



Homo Deus

A Brief History
of Tomorrow

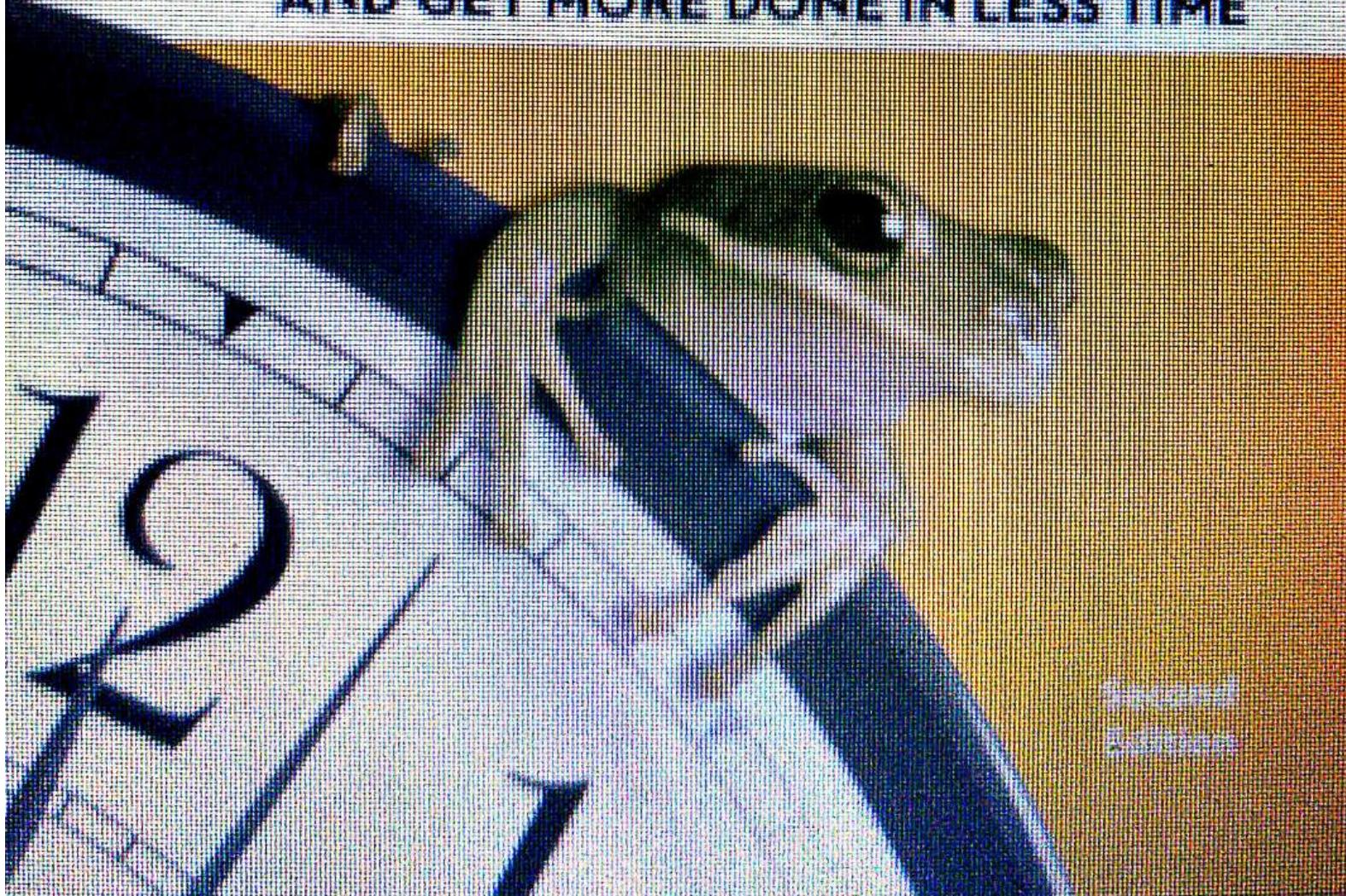
BY YUVAL NOAH HARARI

TRANSLATED BY



EAT THAT FROG!

21 GREAT WAYS TO
STOP PROCRASTINATING
AND GET MORE DONE IN LESS TIME



Second
Edition

BRIAN TRACY

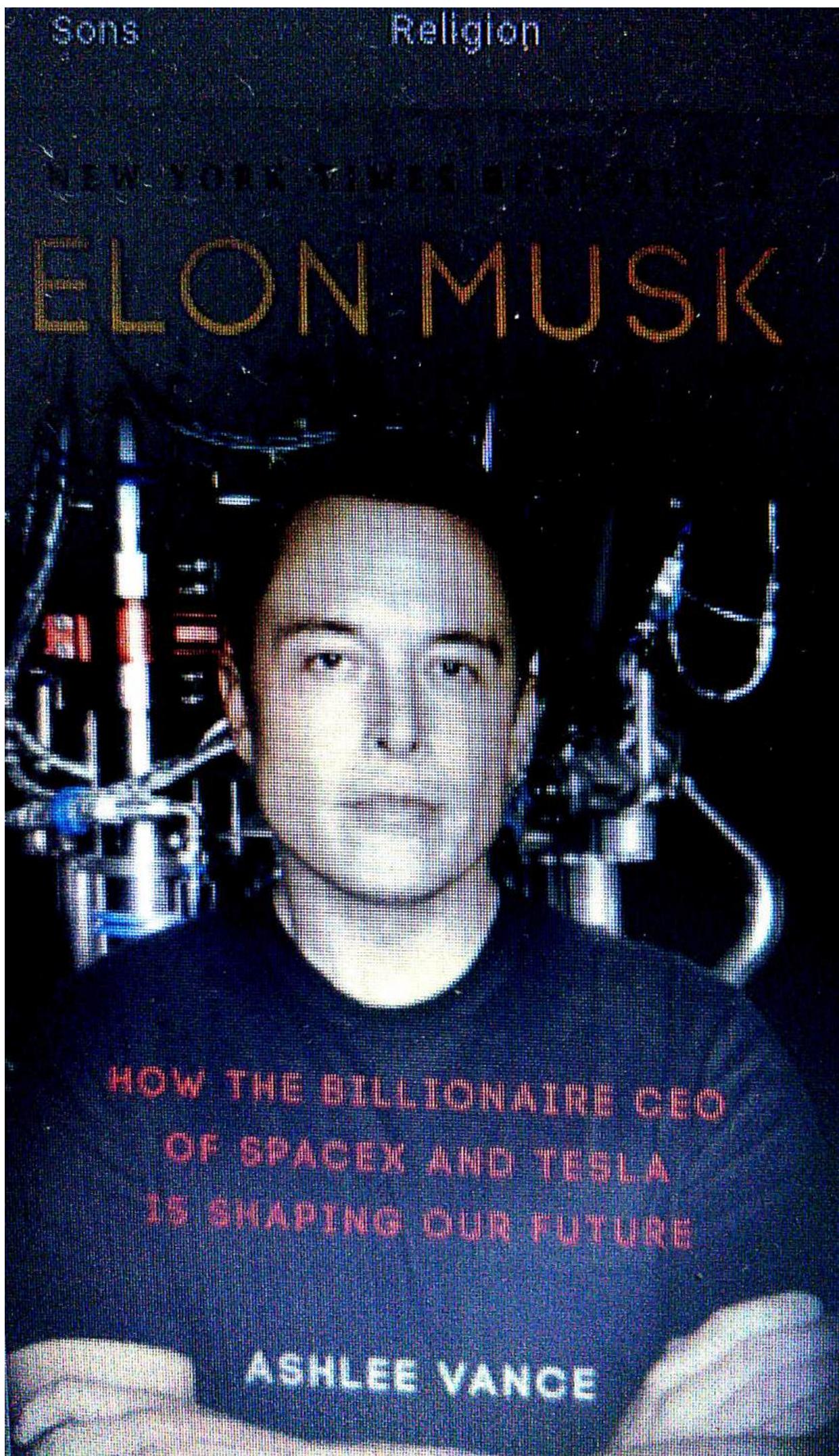
#1 NEW YORK TIMES BESTSELLER

THE SUBTLE ART OF NOT GIVING A F*CK

A COUNTERINTUITIVE APPROACH
TO LIVING A GOOD LIFE

MARK MANSON

THE #1
NEW YORK
TIMES
BESTSELLER
MARK MANSON



Bestselling author of THE 48 LAWS OF POWER

WORLDS GREATEST SELLER

MASSEY

DO Epic Shit

Over
600M views
across social
media

20%
OFF

Ankur Warikoo