

UNIT - I

1

Structure of Words and Documents

Part : I - Finding the Structure of Words

Q.1 Write a short note on natural language processing.

OR Discuss what is NLP.

▣ **Ans. :**

- Natural Language Processing (NLP) is a cross disciplinary field of linguistics, computer science, and artificial intelligence. It is related to the interactions between digitalized computing devices and human language or precisely natural language.
- The field of natural language processing deals with designing and programming digital computational devices (particularly computers) to process and analyse large amounts of natural language data.
- Natural languages take different forms, such as writing, speech or signing. They are unrestricted from constructed and formal languages such as those used to program computers or to study logic. As a result of its natural language data is highly unstructured in nature.
- For example, NLP makes it possible for computers to read text, hear speech, interpret it, measure sentiment and determine which parts are important.

Q.2 Discuss the spectrum of natural languages.

OR Why understanding natural language is challenging

▣ **Ans. :**

- The spectrum of natural languages is very wide. As per the linguistic science There are thousands of spoken languages in the world. These languages can be grouped together as members of a language family.
- There are three main language in the world
 - Indo-European (Includes English)
 - Sino-Tibetan (Includes Chinese)
 - Afro-Asiatic (Includes Arabic)
- Natural language is a complicated thing. In many languages, words are delimited in the orthography by whitespace and punctuation.
- Some languages use word forms that need not change much with the changing context. On the other hand, there are languages that are highly sensitive about the choice of word forms according to context.

- For some of the languages, the context does not impact the gender of the noun, while some languages do not have the concept of gender.
- Natural languages show structure (namely grammar) of different kinds and complexity. It consists of more elementary components whose co-occurrence in context refines the purpose when used in isolation, but they extend it further to meaningful relations between other components in the sentence.
- As a result, understanding natural language in word blocks is not a viable approach. However, the first level understanding of a word is very important.

Q.3 Write a note on word morphology in natural language.

OR Discuss why it is important to understand words in natural language.

▣ **Ans. :** Words are the most indicative blocks of a natural sentence. However, they are tricky to define. This is primarily due to ambiguity and contextual meaning of words in sentences. Knowing how to work with words allows, the development of syntactic and semantic understanding.

- The process of understanding words in any natural language involves morphology, word structure and its linguistic expression.
 - Morphology is the study the variable forms and functions of words, while syntax is concerned with the arrangement of words into phrases, clauses, and sentences.
 - Word structure constraints due to pronunciation are described by phonology, whereas conventions for writing constitute the orthography of a language.
 - The linguistic expression is its semantics, and etymology and lexicology cover especially the evolution of words and explain the semantic, morphological and other links among them.
- Word structure constraints due to pronunciation are described by phonology, whereas conventions for writing constitute the orthography of a language. In this unit we would discuss morphological parsing.
- The technique of discovering of word structure is called morphological parsing. It is used to
 - Identify words of distinct types in natural languages.
 - Model different internal structure of words in connection with the grammatical properties and lexical concepts.

1.1 Words and their Components

Q.4 Discuss the components of words from NLP perspective.

▣ **Ans. :**

- Words are defined as the smallest linguistic units in a natural language that can form a complete utterance by themselves.
- The minimal parts of words that deliver aspects of meaning to them are called morphemes.
- Depending on the means of communication, morphemes are
 - Spelled out via **graphemes** (symbols of writing) such as letters or characters - or
 - Realized through **phonemes** (the distinctive units of sound in spoken language).

- Specific to a particular language the exact boundaries of separating words from morphemes and phrases is varied and it is specific to that language. Here is an example with nouns as a valid word in English used to understand the above concept. Refer Fig. Q.4.1.

Noun	Noun + s (plural)	Noun + s (possessive)	Pronunciation(both)
thrush	thrushes	thrush's	iz
toy	toys	toy's	z
block	blocks	block's	s

Fig. Q.4.1 : Example to separate words from morphemes

Q.5 Write a short note on tokens.

▣ **Ans. :** Tokens are syntactic words. Let's consider a simple sentence given below.

I don't want to buy this product.

- In the above sentence for reasons of generality, linguists prefer to analyse *don't* as two syntactic words (do not), or tokens, each of which has its independent role and can be reverted to its normalized form. On the other hand, all other words in this sentence are treated as an independent single token.
- In English, such tokenization and normalization may be applied to a limited set of cases. However, in other languages, these phenomena have to be treated in a less trivial manner.

Q.6 Discuss clitics in morphology.

OR Explain the concept of clitics from morphology perspective.

▣ **Ans. :** In morphology and syntax, a clitic is a morpheme that has syntactic characteristics of a word, but depends phonologically on another word or phrase.

- It is syntactically independent but phonologically dependent - always attached to a host.
- A clitic is pronounced like an affix, but plays a syntactic role at the phrase level.
- In other words, clitics have the form of affixes, but the distribution of function words.
- For example, the contracted forms of the auxiliary verbs in I'm and we've are clitics.
- As a result, underlying lexical or syntactic units are thereby transformed into one compact string of letters and no longer appear as distinct words.

▣ **Example-**

In Sanskrit language the token Namaste is a Clitic formed due to Sandhi of two tokens < namah + te > (देवनागरीनमः + ते = नमस्ते).

The sandhi changed token

namah → *namas*(नमः → नमस्).

- Tokens behaving in this way can be found in various languages like Latin, Ancient Greek, Chinese, Japanese, Sanskrit, Czech, Tamil, Telugu, Estonian and many more.
- In such languages, tokenization, also known as word segmentation, is the fundamental step of morphological analysis and a prerequisite for most language processing applications.

Q.7 Explain the importance of lexemes as a linguistic form.

OR Discuss the lexemes.

▣ **Ans. :**

- In a natural language a word often denotes one linguistic form in the given context and the concept behind the form and the set of alternative forms that can express it. Such sets are called lexemes or lexical items. They together form the lexicons of a language.
- Lexemes can be divided by their behaviour into the lexical categories of verbs, nouns, adjectives, conjunctions, particles, or other parts of speech.
- The citation form of a lexeme is also called its lemma.
- The notion of the lexeme is central to morphology and it is the basis for defining other concepts in the morphology. For example, the difference between inflection and derivation can be stated in terms of lexemes :
 - Inflectional rules relate a lexeme to its forms.
 - When we convert a word into its other forms, such as turning the singular mouse into the plural mice or mouses, we say we inflect the lexeme.
 - Derivational rules relate a lexeme to another lexeme.
 - When we transform a lexeme into another one that is morphologically related, regardless of its lexical category, for instance, the nouns consumer and consumption are derived from the verb to consume.

Q.8 Discuss morphemes as the smallest meaningful unit in a language.

OR Explain the difference between morphemes and words in NLP.

▣ **Ans. :**

- There are different opinions on whether and how to associate the properties of word forms with their structural components. These components are usually called segments or morphs.
- The morphs that by themselves represent some aspect of the meaning of a word are called morphemes of some function.
- A morpheme is the smallest meaningful unit in a language.
- A morpheme is not identical to a word. The main difference between them is that a morpheme sometimes does not stand alone, but a word, by definition, always stands alone.
- When a morpheme stands by itself, it is considered as a root because it has a meaning of its own (such as the morpheme dog). When it depends on another morpheme to express an idea, it is an affix because it has a grammatical function (such as the -s in dogs to indicate that it is plural)
- Natural languages use different techniques by which morphs and morphemes are combined into word forms. The simplest morphological process concatenates morphs one by one.
- For example, as in the word, mis-manage-ment-s, where manage is a free lexical morpheme and the other elements are morphemes adding some meaning to the whole word.
- For example, in Korean language, many morphemes change their forms with the phonological context. Fig. Q.8.1 shows some Korean morphemes -ess-, -ass-, -yess- indicating past tense.

	concatenated		contracted		
(a)	보았-	po-ass-	봤-	pwass-	'have seen'
(b)	가지었-	ka.ci-ess-	가졌-	ka.cyess-	'have taken'
(c)	하였-	ha-yess-	했-	hayss-	'have done'
(d)	되었-	toy-ess-	됐-	twayss-	'have become'
(e)	놓았-	noh-ass-	놀-	nwass-	'have put'

Fig. Q.8.1 : Korean morphemes indicating past tense

 Allomorphs

- The alternative forms of a morpheme are termed allomorphs.
- Allomorphs are variants of a morpheme that differ in pronunciation but are semantically identical. For example, the English plural marker *-(e)s* of regular nouns can be pronounced /-s/ (bags), (bushes), depending on the final sound of the noun's plural form.

Q.9 Write a short note on following terminologies :

- | | |
|-------------------------|-------------------------------------|
| (a) Typology | (b) Isolating, or analytic typology |
| (c) Synthetic languages | (d) Agglutinative languages |
| (e) Fusional languages | (f) Nonlinear languages |

 Ans. : (a) Typology

- Typology (or Morphological typology) is a way of classifying the languages in the world. It groups languages according to their common morphological structures. Typology organizes languages on the basis of how those languages form words by combining morphemes.
- The typology that is based on quantitative relations between words, their morphemes, and their features as follows.

 (b) Isolating, or analytic typology

- These languages include no or relatively few words that has more than one morpheme. Examples are Chinese, Vietnamese, and Thai.
- Analytic languages show a low ratio of morphemes to words, nearly one-to-one.
- Sentences in analytic languages are composed of independent root morphemes.
- Grammatical relations between words are expressed by separate words where they might otherwise be expressed by affixes, which are present to a minimal degree in such languages.
- Some analytic tendencies are also found in languages like English and Afrikaans.

 (c) Synthetic languages

- These can combine more morphemes in one word and are further divided into agglutinative and fusional languages.
- The morphemes may be distinguishable from the root, or they may not. They may be fused with it or among themselves.

- Word order is less important for these languages than it is for analytic languages, since individual words express the grammatical relations that would otherwise be indicated by syntax.
- In addition, there tends to be a high degree agreement or cross-reference between different parts of the sentence.
- Therefore, morphology in synthetic languages is more important than syntax.
- Most Indo-European languages are moderately synthetic.

□ (d) Agglutinative languages

- These languages have morphemes associated with only a single function at a time.
- Agglutinative languages have words containing several morphemes that are always clearly differentiable from one another.
- Each morpheme represents only one grammatical meaning and the boundaries between those morphemes are easily demarcated.
- The bound morphemes are affixes and they may be individually identified.
- Agglutinative languages tend to have a high number of morphemes per word, and their morphology is usually highly regular.
- Agglutinative languages include Finnish, Hungarian, Turkish, Mongolian, Korean, Japanese, Indonesian, Tamil etc.

□ (e) Fusional languages

- These languages are defined by their feature-per-morpheme ratio higher than other languages.
- Morphemes in fusional languages are not readily distinguishable from the root or among themselves.
- Several grammatical bits of meaning may be fused into one affix.
- Morphemes may also be expressed by internal phonological changes in the root.
- The Indo-European and Semitic languages are the most typically cited examples of fusional languages.
- Examples of fusional Indo-European languages are : Kashmiri, Sanskrit, Pashto, New Indo-Aryan languages such as Punjabi, Hindustani, Bengali; Greek (classical and modern), Latin, Italian, French, Spanish, Portuguese, Romanian, Irish, German, Faroese, Icelandic, Albanian and all Balto-Slavic languages.

Concatenative languages

- These languages link morphs and morphemes one after another.

□ (f) Nonlinear languages

- Nonlinear languages allow structural components to merge non-sequentially to apply tonal morphemes or change the consonantal or vocalic templates of words.
- It is also called discontinuous morphology and intoflection, is a form of word formation and inflection in which the root is modified and which does not involve stringing morphemes together sequentially.

- For example, in English, mostly plurals are usually formed by adding the suffix - s, certain words use nonconcatenative processes for their plural forms as

foot → feet

- Many irregular verbs form their past tenses, past participles or both in the same manner :

freeze → froze → frozen

- This specific form of nonconcatenative morphology is known as base modification or ablaut, a form in which part of the root undergoes a phonological change without necessarily adding new phonological material

- For example the English stem song, results in the four distinct words as

Sing → sang → sung → song

1.2 Issues and Challenges

Q.10 Explain the importance of morphological parsing and modelling in NLP.

▣ Ans. :

- Morphological parsing helps to eliminate or improve the inconsistency of word forms. It is required to provide higher-level linguistic units whose lexical and morphological properties are explicit and well defined.
- Every Natural language inherently has some irregularity and ambiguity. Morphological parsing attempts to remove unnecessary irregularity and control ambiguity.
- **Irregularities**
 - In this context irregularity means existence of such forms and structures that are not described appropriately by a prototypical linguistic model.
 - Some irregularities can be understood by redesigning the model and improving its rules, but other lexically dependent irregularities often cannot be generalized.
- **Ambiguity**
 - Ambiguity is an inability in interpretation of expressions of language.
 - Accidental ambiguity and ambiguity due to lexemes with multiple senses, cause syncretism, or systematic ambiguity.
- Morphological modelling also faces the problem of productivity and creativity in language. This gives birth to unconventional but perfectly meaningful new words or new senses to the language.
- Because these newly coined words are not present in the lexical and morphological properties, such words will remain completely unparsed in morphological system. This unknown word problem is particularly severe in speech or writing.
- The morphological modelling is unable to parse a word, that comes from an expected domain of the linguistic model. This happens mostly when special terms or foreign names are involved in the discourse or when multiple languages or dialects are mixed together.

- For example, in English, mostly plurals are usually formed by adding the suffix - s, certain words use nonconcatenative processes for their plural forms as

foot → feet

- Many irregular verbs form their past tenses, past participles or both in the same manner :

freeze → froze → frozen

- This specific form of nonconcatenative morphology is known as base modification or ablaut, a form in which part of the root undergoes a phonological change without necessarily adding new phonological material
- For example the English stem song, results in the four distinct words as

Sing → sang → sung → song

1.2 Issues and Challenges

Q.10 Explain the importance of morphological parsing and modelling in NLP.

▣ Ans. :

- Morphological parsing helps to eliminate or improve the inconsistency of word forms. It is required to provide higher-level linguistic units whose lexical and morphological properties are explicit and well defined.
- Every Natural language inherently has some irregularity and ambiguity. Morphological parsing attempts to remove unnecessary irregularity and control ambiguity.
- **Irregularities**
 - In this context irregularity means existence of such forms and structures that are not described appropriately by a prototypical linguistic model.
 - Some irregularities can be understood by redesigning the model and improving its rules, but other lexically dependent irregularities often cannot be generalized.
- **Ambiguity**
 - Ambiguity is an inability in interpretation of expressions of language.
 - Accidental ambiguity and ambiguity due to lexemes with multiple senses, cause syncretism, or systematic ambiguity.
- Morphological modelling also faces the problem of productivity and creativity in language. This gives birth to unconventional but perfectly meaningful new words or new senses to the language.
- Because these newly coined words are not present in the lexical and morphological properties, such words will remain completely unparsed in morphological system. This unknown word problem is particularly severe in speech or writing.
- The morphological modelling is unable to parse a word, that comes from an expected domain of the linguistic model. This happens mostly when special terms or foreign names are involved in the discourse or when multiple languages or dialects are mixed together.

Q.11 Explain morphological irregularities in NLP.

OR Discuss how the morphological irregularities are removed.

▣ **Ans. :**

- The design principles of the morphological model are very important to control the irregularities in words.
- Morphological parsing is designed for generalization and abstraction of words to make the model simple and yet powerful.
- However, the immediate descriptions of given for a word may not be the final ones, due to
 - Inadequate accuracy description
 - Inappropriate complexity of morphological model
 - Need of improved formulations
- **Removal of morphological irregularities**
 - A deeper study of the morphological processes is essential for mastering the whole morphological and phonological system.
 - Morphophonemic templates capture morphological processes. It is done by organizing stem patterns and generic affixes.
 - These templates are designed without any context-dependent variation of the affixes or ad hoc modification of the stems.
 - A very terse merge rules ensure that morphophonemic templates can be converted into exactly the surface forms namely, orthographic and phonological.
 - Applying the merge rules is independent of and irrespective of any grammatical parameters or information other than that contained in a template.
 - Thus, most morphological irregularities in the morphophonemic templates are successfully removed.

Q.12 Discuss morphological irregularities in any two natural languages.

▣ **Ans. :**

- **Morphological irregularities in Arabic**
 - Morphophonemic templates can be used for discovering the regularity of Arabic morphology where uniform structural operations apply to different kinds of stems.
 - Some irregularities are bound to particular lexemes or contexts, and cannot be accounted for by general rules.
- **Morphological irregularities in Korean**
 - Korean irregular verbs provide examples of such irregularities. Korean shows exceptional constraints on the selection of grammatical morphemes.
 - Korean language features lexically dependent stem alternation.

• Morphological irregularities in other Natural languages

- It is hard to find irregular inflection in agglutinative languages : Two irregular verbs in Japanese, one in Finnish.
- These languages are abundant with morphological alternations that are formalized by precise phonological rules.

Q.13 What is morphological ambiguity ? Discuss at least two examples.**▣ Ans. :**

- Morphological ambiguity is the possibility that word forms be understood in multiple ways out of the context.
- Words forms that look the same but have distinct functions or meaning are called homonyms.
- Ambiguity is present in all aspects of morphological processing and language processing at large.
- Morphological parsing cannot complete disambiguation of words in their context, but it can control the valid interpretations of a given word form.
- **Morphological Ambiguity in Korean**
 - In Korean, homonyms are one of the most problematic objects in morphological analysis. This is because they prevail all around frequent lexical items.
- **Morphological Ambiguity in Arabic**
 - Arabic has rich derivational and inflectional morphology. Because Arabic script usually does not encode short vowels and omits yet diacritical marks, its morphological ambiguity is considerably increased. In addition, Arabic orthography collapses certain word forms together.
 - The problem of morphological disambiguation of Arabic encompasses
 - ⇒ The resolution of the structural components of words
 - ⇒ Actual morphosyntactic properties
 - ⇒ Tokenization and normalization
 - ⇒ Lemmatization, stemming
 - ⇒ Diacritization
- **Morphological ambiguity in Sanskrit**
 - When inflected syntactic words are combined in an utterance, additional phonological and orthographic changes can take place.
 - In Sanskrit, one such euphony rule is known as external sandhi. Inverting sandhi during tokenization is usually nondeterministic as it can provide multiple solutions.
- In any language, tokenization decisions may impose constraints on the morphosyntactic properties of the tokens being reconstructed.
- The morphological phenomenon that some words or word classes show instances of systematic homonymy is called syncretism. In particular, homonymy can occur due to neutralization and unaffectedness of words.

Q.14 What is morphological productivity ?

OR Discuss the morphological productivity.

OR Discuss the competence versus performance duality by noam chomsky in the context of morphological productivity.

▣ **Ans. :**

- In a natural language as a system (langue), structural devices like recursion, iteration, or compounding allow to produce an infinite set of concrete linguistic utterances.
- This general potential holds for morphological processes as well and is called morphological productivity.
- In a perspective natural language can be seen as a collection of utterances (parole) pronounced or written (performance). Hence for the linguistic corpora, parole and performance data set is practical.
- Such corpora are a finite collection of linguistic data that are studied with empirical methods. It can be used for comparison when linguistic models are developed.

Q.15 Discuss "80/20 rule," of linguistic word corpus ?

OR Write a note on "80/20 rule" of linguistic word corpus.

▣ **Ans. :**

- Linguistic corpora are a finite collection of linguistic data that are studied with empirical methods.
- The set of word forms found in the corpus of a language is referred as its vocabulary.
- The members of this set are word types, whereas every original instance of a word form is a word token.
- The distribution these words or other elements of language follows the "80/20 rule," also known as the law of the vital few.
- It says that most of the word tokens in a given corpus can be identified with just a couple of word types in its vocabulary, and words from the rest of the vocabulary occur much less or rarely in the corpus.
- New, unexpected words will always appear in the linguistic data only when it is expanded or enlarged.

Q.16 Discuss how creativity and the issue of unknown words meet to enhance the morphological productivity in a natural language.

OR Discuss how the newly coined word google has enhanced the morphological productivity of many natural languages.

OR Discuss unexpected words will always appear in the linguistic data only when it is expanded or enlarged.

▣ **Ans. :**

- The word googol is a dictionary word in English. It means something that is a made-up word denoting the number "one followed by one hundred zeros,".
- The name of the company Google is an inadvertent misspelling thereof. However, both of these words successfully entered the lexicon of English where morphological productivity started working.
- Today we understand English verb to google and nouns like googling or even googlish or googleology.

- This new word google is adopted by other languages, too. This has triggered their own morphological processes.
- In Czech, one says googlovat, googlit 'to google' or vygooglovat, vygooglit 'to google out', googlování 'googling', and so on.
- In Arabic, the names are transcribed as ġūġūl 'googol' and ġūġil 'Google'.
- Thus we can observe that unexpected words in a language will always appear in the linguistic data only when it is expanded or enlarged.

1.3 Morphological Models

Q.17 Discuss the motivation of using domain-specific languages.

OR What is Domain Specific Language (DSL) ?

▣ Ans. :

- A Domain Specific Language (DSL) is a specialized programming language that is used for a single purpose.
- Various domain-specific languages have been created for achieving intuitive and minimal programming effort.
- Pragmatically, a DSL may be specialized to a particular problem domain, a particular problem representation technique, a particular solution technique, or other aspects of a domain.
- These special-purpose languages usually introduce idiosyncratic notations of programs and are interpreted using some restricted model of computation.
- The motivation for this approach lies in the fact that, historically, computational resources were too limited compared to the requirements and complexity of the tasks being solved.
- Other motivations are theoretical given that finding a simple, accurate and yet generalizing model for the practical use in the specific domain.
- The design objective of DSL is to get be pure, intuitive, adequate, complete, reusable and elegant language.
- Examples of such domain-specific programming languages are HTML, SQL, AWK, GDL, etc.

Q.18 Why dictionary lookup is considered as one of the effective Morphological model ?

OR Discuss dictionary as a morphological model.

▣ Ans. :

- Morphological model needs a system in which analysing a word form is reduced kept in sync with more sophisticated models of the language. Dictionaries, Databases and Lists are some of such forms.
- A dictionary is understood as a data structure that directly enables obtaining some precomputed results i.e. word analyses.
- The data structure can be optimized for efficient lookup, and the results can be shared.

- Lookup operations with dictionaries are relatively simple and usually quick. Dictionaries can be implemented, for instance, as lists, binary search trees, tries, hash tables, etc.
- Hence dictionary lookup is considered as one of the effective Morphological models.

Q.19 What are the drawbacks of enumerative Morphological model ?

▣ Ans. :

- Enumerative list is a set of associations between word forms and their desired descriptions.
- It is declared by plain enumeration. Hence the coverage of the model is finite and the generative potential of the language is not exploited.
- Development, lookup and verification of the association list is tedious, liable to errors, inefficient and inaccurate unless the data are retrieved automatically from large and reliable linguistic resources.
- Despite all that, an enumerative model is often sufficient for the given purpose, deals easily with exceptions, and can implement even complex morphology.

Q.20 Write a short note on finite-state morphological.

▣ Ans. :

- Finite-state morphological models are the morphological models in which the specifications written by human programmers are directly compiled into finite-state transducers.
- The finite state morphological models can be used for multiple natural languages.
- The two popular online tools supporting this approach are XFST (Xerox Finite-State Tool) and LexTools.

Q.21 Discuss finite state transducers.

OR Discuss how the finite state transducers can translate the infinite regular language.

▣ Ans. :

- Finite-state transducers are computational devices extending the power of finite-state automata.
- They consist of a finite set of nodes connected by directed edges labeled with pairs of input and output symbols.
- In such a network or graph, nodes are also called states, while edges are called arcs.
- Traversing the network from the set of initial states to the set of final states along the arcs is equivalent to reading the sequences of encountered input symbols and writing the sequences of corresponding output symbols.
- The set of possible sequences accepted by the transducer defines the input language; the set of possible sequences emitted by the transducer defines the output language.
- For example, a finite-state transducer could translate the infinite regular language consisting of the Sanskrit words *pita*, *prapita*, *praprapita*, ... to the matching words in the infinite regular English language words defined as *father*, *grand-father*, *great-grand-father*.
- In finite-state transducers it is possible to invert the domain and the range of a relation, that is, exchange the input and the output.
- In finite-state computational morphology, it is common to refer to the input word forms as surface strings and to the output descriptions as lexical strings.

Part : II - Finding Structure of Document


Q.22 What is a structure of document ?

OR What is a document structure.

OR Write a short note of document structure.

▣ **Ans. :**

- In human language, words and sentences do not appear randomly but usually have a structure.
- For example, combinations of words form sentences - meaningful grammatical units, such as statements, requests, and commands.
- Likewise, in written text, sentences form paragraphs - self-contained units of discourse about a particular point or idea.

 **1.4 Introduction**

Q.23 Discuss the importance of document structure in human language.

OR Why document structure is important in NLP ?

▣ **Ans. :**

- In human language or natural language, words and sentences usually have a structure. This can be combinations of words form sentences - meaningful grammatical units, such as statements, requests, and commands.
- Similarly, in written text, paragraphs are the self-contained units about a point or an idea, which is expressed in the form of group of sentences. Following are the some of the reasons why document structure is important in human languages and therefore for natural language processing.
- When the structure of documents is extracted, it makes the further processing of text easy in NLP. The NLP tasks that depends on the document structure are, parsing, machine translation and semantic role labelling in sentences.
- To improve the reliability of Automatic Speech Recognition (ASR) and human readability, it is important to identify the sentence boundary annotation. Document structure helps in this process.
- Document structure helps in breaking apart the input text or speech into topically coherent blocks that provides better organization and indexing of the data.
- Thus, in most speech and language processing applications extracting the structure of textual and audio documents is a meaningful and necessary pre-step.

Q.24 Write a note on sentence boundary detection.

OR What is sentence boundary detection ?

▣ **Ans. :**

- Sentence boundary detection is the problem in natural language processing of deciding where sentences begin and end.

- Sentence detection is an important task, which should be performed at the beginning of a text processing pipeline.
- Sentence boundary detection (also called sentence segmentation) deals with automatically segmenting a sequence of word tokens into sentence units.
- Natural language processing tools often require their input to be divided into sentences; however, sentence boundary identification can be challenging due to the potential ambiguity of punctuation marks.
- In written text in English and some other languages, the beginning of a sentence is usually marked with an uppercase letter, and the end of a sentence is explicitly marked with a period (.), a question mark (?), an exclamation mark or another type of punctuation.
- However, in addition to their role as sentence boundary markers, capitalized initial letters are used to distinguish proper nouns, periods are used in abbreviations and numbers and other punctuation marks are used inside proper names.
- A character-wise analysis of text allows for a distinction between period characters that are enclosed between two alphanumeric characters, and period characters that are followed by at least one, non-alphabetic character, such as a further punctuation sign, a space, tab or new line.
- There are various challenges associated with SBD, for written as well as spoken text and code switching.

Q.25 Discuss the challenges of sentence boundary detection in written text.

▣ **Ans. :**

- Ambiguous abbreviations and capitalizations are the most common problems of sentence segmentation in written text.
- Quoted sentences are more complex and problematic. The primary reason for this is the speaker may have uttered multiple sentences and sentence boundaries inside the quotes are also marked with punctuation marks.
- As a result of this an automatic method of sentence boundary detection may result in cutting some sentences incorrectly. In case if the preceding sentence is spoken instead of written, prosodic cues usually mark structure.
- “Spontaneously” written texts, such as Short Message Service (SMS) texts or Instant Messaging (IM) texts, tend to be nongrammatical and have poorly used or missing punctuation, which makes sentence segmentation even more challenging.
- The automatic systems, such as optical character recognition (OCR) or ASR, aims to translate images of handwritten, typewritten, or printed text or spoken utterances into machine-editable text.
- When the sentences comes from such automatic system, the finding of sentence boundaries must deal with the errors of these systems as well.
- For example, OCR system easily confuses periods and commas and can result in meaningless sentences. ASR transcripts typically lack punctuation marks and are usually mono-case.

Q.26 Discuss the challenges of sentence boundary detection in spoken / conversational text.

▣ **Ans. :**

- For conversational speech or text or multiparty meetings with ungrammatical sentences and disfluencies, in most cases it is not clear where the boundaries are.

- The problem may be redefined for the conversational domain as the task of dialog act segmentation. This is because dialog acts are better defined for conversational speech using a number of mark-up standards such as Dialog Act Mark-up in Several Layers (DAMSL).
- For example, the sentence I think so but you should also ask him may be a grammatical sentence as a whole, but for DAMSL and MRDA standards, there are two dialog act tags, one affirmation and one suggestion. Such a modification may be needed for conversation analysis, such as speaker role detection or sentiment analysis. This task can be seen as a semantic boundary detection task instead of syntactic.

Q.27 What is code switching ? Why it is considered as a problem in sentence boundary detection ?

▣ **Ans. :**

- Code switching - that is, the use of words, phrases, or sentences from multiple languages by multilingual speakers - is another problem that can affect the characteristics of sentences. For example, when switching to a different language, the writer can either keep the punctuation rules from the first language or resort to the code of the second language (e.g., Spanish uses the inverted question mark to precede questions).
- Code switching also affects technical texts for which the meanings of punctuation signs can be redefined, as in Uniform Resource Locators (URLs), programming languages, and mathematics. We must detect and parse those specific constructs in order to process technical texts adequately.
- Conventional rule-based sentence segmentation systems in well-formed texts rely on patterns to identify potential ends of sentences and lists of abbreviations for disambiguating them.
- Although rules cover most of these cases, they do not address unknown abbreviations, abbreviations at the ends of sentences, or typos in the input text.
- Furthermore, such rules are not robust to text that is not well formed, such as forums, chats, and blogs, or to spoken input that completely lacks typographic cues. Moreover, each language requires a specific set of rules.
- Hence code switching is considered as a problem in sentence boundary detection.

Q.28 How sentence segmentation as a classification problem is more effective than a rule based problem ?

▣ **Ans. :**

- Conventional rule-based sentence segmentation systems in well-formed texts rely on patterns to identify potential ends of sentences and lists of abbreviations for disambiguating them.
- Sentence segmentation in text usually uses the punctuation marks as delimiters and aims to categorize them as sentence ending/beginning or not. On the other hand, for speech input, all word boundaries are usually considered as candidate sentence boundaries.
- Although rules cover most of these cases, they do not address unknown abbreviations, abbreviations at the ends of sentences, or typos in the input text.
- Furthermore, such rules are not robust to text that is not well formed, such as forums, chats, and blogs, or to spoken input that completely lacks typographic cues. Moreover, each language requires a specific set of rules.
- To improve on such a rule-based approach, sentence segmentation is stated as a classification problem. Given training data where all sentence boundaries are marked, we can train a classifier to recognize them.

Q.29 What is topic boundary segmentation ?

OR How automatic topic boundary segmentation works ?

▣ **Ans. :**

- Topic segmentation (sometimes called discourse or text segmentation) is the task of automatically dividing a stream of text or speech into topically homogeneous blocks.
- That is, given a sequence of (written or spoken) words, the aim of topic segmentation is to find the boundaries where topics.
- Topic segmentation is an important task for various language-understanding applications, such as information extraction and retrieval and text summarization.
- In information retrieval, if long documents can be segmented into shorter, topically coherent segments, then only the segment that is about the user's query could be retrieved.
- For multiparty meetings, the task of topic segmentation is inspired by discourse analysis.
- For official and well-structured meetings, the topics are segmented according to the agenda items, whereas for more casual conversational-style meetings, the boundaries are less clear.
- For conversational speech, the topic boundaries may not be absolute. Hence they are more complex.
- In text, topic boundaries are usually marked with distinct segmentation cues, such as headlines and paragraph breaks. These cues are absent in speech. However, speech provides other cues, such as pause duration and speaker changes.
- Topic segmentation is a nontrivial problem without a very high human agreement because of many natural-language-related issues and hence requires a good definition of topic categories and their granularities.

1.5 Methods

Q.30 Discuss Sentence / Topic segmentation as a boundary classification problem.

OR Why Sentence / Topic segmentation is considered as a boundary classification problem in NLP ?

▣ **Ans. :**

- Sentence segmentation and topic segmentation have mainly been considered as a boundary classification problem.
- For given a boundary candidate (between two-word tokens for sentence segmentation and between two sentences for topic segmentation), the goal is to predict whether or not the candidate is an actual boundary (sentence or topic boundary).
- Formally, let $x \in X$ be the vector of features (the observation) associated with a candidate and $y \in Y$ be the label predicted for that candidate. The label y can be b for boundary and \bar{b} for nonboundary.
- This results in a classification problem : given a set of training examples $\{x, y\}_{\text{train}}$, find a function that will assign the most accurate possible label y of unseen examples x_{unseen} .
- Alternatively to the binary classification problem, it is possible to model boundary types using finer-grained categories.

- Gillick suggested that sentence segmentation in text be framed as a three-class problem : sentence boundary with an abbreviation b^a , without an abbreviation $b^{\bar{a}}$, and abbreviation not at a boundary $b^{\bar{a}}$.
- Similarly, in spoken language, a three-way classification can be made between non-boundaries \bar{b} , statement b^s , and question boundaries b^q .

Q.31 Discuss the method of classification in sentence or topic segmentation.

OR Explain the classification method used in sentence or topic segmentation.

▣ **Ans. :**

- For sentence or topic segmentation, the problem is defined as finding the most probable sentence or topic boundaries.
- The natural unit of sentence segmentation is words and of topic segmentation is sentences, with assumption that assume topics typically do not change in the middle of a sentence.
- The words or sentences are then grouped into contiguous stretches belonging to one sentence or topic - that is, the word or sentence boundaries are classified into sentence or topic boundaries and non-boundaries.
- The classification can be done at each potential boundary i (local modelling); then, the aim is to estimate the most probable boundary type, \hat{y}_i , for each candidate example, x_i :

$$\hat{y}_i = \operatorname{argmax}_{y_i \text{ in } y} P(y_i | x_i)$$

- Here, the $\hat{}$ is used to denote estimated categories, and a variable without a $\hat{}$ is used to show possible categories.
- In local modelling, features can be extracted from the surrounding example context of the candidate boundary to model such dependencies. It is also possible to see the candidate boundaries as a sequence and search for the sequence of boundary types, $\hat{Y} = \hat{y}_1, \dots, \hat{y}_n$, that have the maximum probability given the candidate examples, $X = x_1, \dots, x_n$:

$$\hat{Y} = \operatorname{argmax}_Y P(Y | X)$$

Q.32 Discuss the categorization of methods according to the type of the machine learning algorithm.

OR What are the generative and discriminative categorization methods ?

OR Compare between generative and discriminative categorization methods.

▣ **Ans. :**

▣ **Generative sequence model**

- It estimate the joint distribution of the observations, $P(X, Y)$ (e.g., words, punctuation) and the labels (sentence boundary, topic boundary).
- It requires specific assumptions (such as backoff to account for unseen events) and have good generalization properties.

▣ **Discriminative sequence model**

- It focus on features that characterize the differences between the labeling of the examples.
- Such methods (as described in the following sections) can be used for sentence and topic segmentation in both written and spoken language, with one difference.

- c) In text, the category of all boundaries that do not include a potential end-of-sentence delimiter (period, question mark, exclamation mark) is preset to nonsentence or nontopic,
- d) A category is estimated for only those word boundaries that include a delimiter, whereas in speech, all boundaries between consecutive tokens are usually considered.

Q.33 Explain generative sequence classification methods for sentence and topic segmentation.

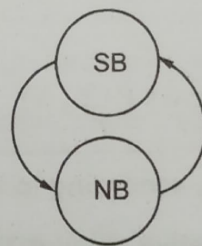
▣ **Ans. :**

- The most commonly used generative sequence classification method for topic and sentence segmentation is the hidden Markov model (HMM).

- The probability is written as the following, using the Bayes rule :

$$\hat{Y} = \operatorname{argmax}_Y P(Y|X) = \operatorname{argmax}_Y \frac{P(X|Y)P(Y)}{P(X)} = \operatorname{argmax}_Y P(X|Y)P(Y)$$

- $P(X)$ in the denominator is dropped because it is fixed for different Y and hence does not change the argument of max.
- The bigram case is modelled by a fully connected m -state Markov model, where m is the number of boundary categories.
- The states emit words (sentences or paragraphs) for sentence (topic) segmentation, and the state sequence that most likely generated the word (sentence) sequence is estimated.
- State transition probabilities, $P(y_i|y_{i-1})$, and state observation likelihoods, $P(x_i|y_i)$, are estimated using the training data.
- The most probable boundary sequence is obtained by dynamic programming.
- Below is the conceptual hidden Markov model for segmentation with two states: one for segment boundaries, one for others.



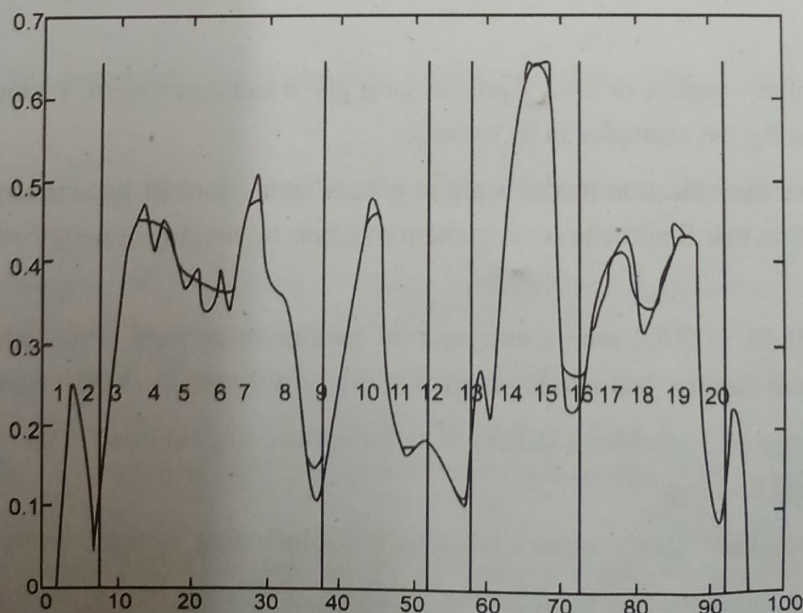
- The bigram case can be extended to higher order n -grams at the cost of an increased complexity.
- For topic segmentation, typically instead of using two states, n states are used, where n is the number of topics. However, it is not possible in HMM to use any information beyond words, such as POS tags of the words or prosodic cues, for speech segmentation.
- Two simple extensions have been proposed : Shriberg et al suggested using explicit states to emit the boundary tokens, hence incorporating nonlexical information via combination with other models.
- For topic segmentation, Tur et al. used the same idea and modeled topic-start and topic-final sections explicitly, which helped greatly for broadcast news topic segmentation. The second extension is inspired from factored language models which capture not only words but also morphological, syntactic, and other information. Guz et al proposed using factored HELM (fHELM) for sentence segmentation using POS tags in addition to words.

Q.34 Discuss discriminative local classification methods.**Ans. :**

- A number of discriminative classification approaches, such as support vector machines, boosting, maximum entropy, and regression, are based on very different machine learning algorithms.
- While discriminative approaches have been shown to outperform generative methods in many speech and language processing tasks, training typically requires iterative optimization.
- In discriminative local classification, each boundary is processed separately with local and contextual features.
- No global (i.e., sentence or document wide) optimization is performed, unlike in sequence classification.
- For sentence segmentation, supervised learning methods have primarily been applied to newspaper articles.
- Many classifiers have been tried for the task : regression trees , neural networks , a C4.5 classification tree , maximum entropy classifiers , support vector machines (SVMs), and naive Bayes classifiers .
- Mikheev treated the sentence segmentation problem as a subtask for POS tagging by assigning a tag to punctuation similar to other tokens . For tagging he employed a combination of HMM and maximum entropy approaches.

Q.35 Write a note on TextTiling method for topic segmentation.**OR Discuss how TextTiling method is used for topic segmentation.****OR Explain block comparison and vocabulary introduction methods for topic segmentation.****Ans. :**

- The popular TextTiling method of Hearst for topic segmentation uses a lexical cohesion metric in a word vector space as an indicator of topic similarity.
- TextTiling can be seen as a local classification method with a single feature of similarity.
- Below Fig. Q.35.1 depicts a typical graph of similarity with respect to consecutive segmentation units. The document is chopped when the similarity is below some threshold.

**Fig. Q.35.1 : TextTiling example**

Originally, two methods for computing the similarity scores were proposed for Text Tiling :

- **Block comparison -**

- It compares adjacent blocks of text to see how similar they are according to how many words the adjacent blocks have in common.
- The block size can be variable, not necessarily looking only at the consecutive blocks but instead at a window.
- Given two blocks, b_1 and b_2 , each having k tokens (sentences or paragraphs), the similarity (or topical cohesion) score is computed by the formula given below.

$$\frac{\sum_t \omega_t b_{1t} \cdot \omega_t b_{2t}}{\sqrt{\sum_t \omega_t^2 b_{1t} \sum_t \omega_t^2 b_{2t}}}$$

where ω_t , b is the weight assigned to term t in block b . The weights can be binary or may be computed using other information retrieval-based metrics such as term frequency.

- **Vocabulary introduction -**

- The vocabulary introduction method, assigns a score to a token-sequence gap on the basis of how many new words are seen in the interval in which it is the midpoint.
- Similar to the block comparison formulation, given two consecutive blocks, b_1 and b_2 , of equal number of words, w , the topical cohesion score is computed with the following formula : Where NumNewTerms(b) returns the number of terms in block b , seen for the first time in text.

$$\frac{\text{NumNewTerms}(b_1) + \text{NumNewTerms}(b_2)}{2 \times w}$$

- This method is extended to exploit latent semantic analysis. Instead of simply looking at all words, researchers worked on the transformed lexical space, which has led to improved results because this approach also captures semantic similarities implicitly.

Q.36 Explain discriminative sequence classification methods.

▣ **Ans. :**

- In segmentation tasks, the sentence or topic decision for a given example (word, sentence, paragraph) highly depends on the decision for the examples in its vicinity.
- Discriminative sequence classification methods are in general extensions of local discriminative models with additional decoding stages that find the best assignment of labels by looking at neighbouring decisions to label an example.
- Conditional Random Fields (CRFs) are an extension of maximum entropy, SVM struct is an extension of SVM to handle structured outputs, and maximum margin. Markov networks (M3N) are extensions of HMMs.
- The Margin Infused Relaxed Algorithm (MIRA) is an online learning approach that requires loading of one sequence at a time during training.
- CRFs have been successful for many sequence labelling tasks, including sentence segmentation in speech.
- CRFs are a class of log-linear models for labelling structures. CRFs are trained by finding the λ parameters that maximize the likelihood of the training data, usually with a regularization term to avoid overfitting.

- Gradient, conjugate gradient, or online methods are used for training.
- Dynamic programming (Viterbi decoding) is used to find the most probable assignment of labels at test time or to compute the $Z(\cdot)$ function.

Q.37 Discuss the hybrid approaches for word classification.

▣ **Ans. :**

- Nonsequential discriminative classification algorithms typically ignore the context, which is critical for the segmentation task.
- While we may add context as a feature or simply use CRFs, which inherently consider context, these approaches are suboptimal when dealing with real-valued features, such as pause duration or pitch range. Most earlier studies simply tackled this problem by binning the feature space either manually or automatically
- An alternative is to use a hybrid classification approach, as suggested by Shriberget al. .
- The main idea is to use the posterior probabilities, $P_c(y_i|x_i)$, for each boundary candidate, obtained from the other classifiers, such as boosting or CRF, by simply converting them to state observation likelihoods by dividing to their priors following the well-known Bayes rule as follows

$$\operatorname{argmax}_{y_i} = \frac{P_c(y_i|x_i)}{P(y_i)} = \operatorname{argmax}_{y_i} P(x_i|y_i)$$

- Applying the Viterbi algorithm to the HMM then returns the most likely segmentation. To handle dynamic ranges of state transition probabilities and observation likelihoods, a weighting scheme as is usually described in the literature can be applied.
- Zimmerman et al. compared various discriminative local classification methods, namely boosting, maximum entropy, and decision trees, along with their hybrid versions for sentence segmentation of multilingual speech. He concluded that hybrid approaches are always superior.

Q.38 What are the extensions for global modeling for sentence segmentation ?

OR

How global modeling for sentence segmentation is carried out using extensions ?

▣ **Ans. :**

- Most approaches to sentence segmentation have focused on recognizing boundaries rather than sentences in themselves.
- This has occurred because of the quadratic number of sentence hypotheses that must be assessed in comparison to the number of boundaries.
- To tackle this problem, input is segmented according to likely sentence boundaries established by a local model. Later it is trained as a re-ranker on the n-best lists.
- This approach allows leveraging of sentence-level features such as scores from a syntactic parser or global prosodic features.
- Favre et al. proposed to extend this concept to a pruned sentence lattice, which allows combining local scores with sentence-level scores in a more efficient manner.

1.6 Complexity of Approaches

Q.39 Discuss how the complexity of sentence/topic segmentation is evaluated.

Ans. :

- Sentence/topic segmentation approaches can be rated in terms of complexity (time and memory) of their training and prediction algorithms and in terms of their performance on real-world datasets. Some may also require specific pre-processing, such as converting or normalizing continuous features to discrete features.
- **Discriminative approach**
 - a) In terms of complexity, training of discriminative approaches is more complex than training of generative ones because they require multiple passes over the training data to adjust for their feature weights.
- **Generative models**
 - b) Generative models such as HELMs can handle multiple orders of magnitude larger training sets and benefit, for instance, from decades of news wire transcripts. But they do not cope well with unseen events.
- **Discriminative classifiers**
 - c) They allow for a wider variety of features and perform better on smaller training sets.
 - d) Predicting with discriminative classifiers is also slower, even though the models are relatively simple (linear or log-linear), because it is dominated by the cost of extracting more features.
- **Sequence approaches**
 - e) Compared to local approaches, sequence approaches bring the additional complexity of decoding: finding the best sequence of decisions requires evaluating all possible sequences of decisions.
 - f) Fortunately, conditional independence assumptions allow the use of dynamic programming to trade time for memory and decode in polynomial time.
 - g) This complexity is then exponential in the order of the model (number of boundary candidates processed together) and the number of classes (number of boundary states).
- **Discriminative sequence classifiers,**
 - h) For example CRFs, also need to repeatedly perform inference on the training data, which might become expensive.

1.7 Performance of the Approaches

Q.40 Discuss the performance of sentence segmentation approaches in detail.

OR Write a short note on :

- a) Sentence segmentation in speech
- b) Sentence segmentation in text
- c) Sentence segmentation in speech

Ans. :

a) Sentence segmentation in speech

- For sentence segmentation in speech, performance is usually evaluated using -
 - 1) The error rate (ratio of number of errors to the number of examples)
 - 2) F1-measure (the harmonic mean of recall and precision)

where

- 1) Recall is defined as the ratio of the number of correctly returned sentence boundaries to the number of sentence boundaries in the reference annotations.
- 2) Precision is the ratio of the number of correctly returned sentence boundaries to the number of all automatically estimated sentence boundaries), and the National Institute of Standards and Technology (NIST) error rate (number of candidates wrongly labeled divided by the number of actual boundaries).

b) **Sentence segmentation in text**

- For sentence segmentation in text, researchers have reported error rate results on a subset of the Wall Street Journal Corpus of about 27,000 sentences.
- For instance, Mikheev reports that his rule-based system performs at an error rate of 1.41 %.
- The addition of an abbreviation list to this system lowers its error rate to 0.45 % and combining it with a supervised classifier using POS tag features leads to an error rate of 0.31 %.
- Without requiring handcrafted rules or an abbreviation list, Gillick's SVM-based system obtains even fewer errors, at 0.25 %.
- Even though the error rates presented seem low, sentence segmentation is one of the first processing steps for any NLP task, and each error impacts subsequent steps, especially if the resulting sentences are presented to the user as for example, in extractive summarization.

c) **Sentence segmentation in speech**

- For sentence segmentation in speech, Doss et al. report on the Mandarin TDT4 Multilingual Broadcast News Speech Corpus, an F1-measure using the same set of features is as of
 - 69.1 % for a MaxEnt classifier
 - 72.6 % with Adaboost
 - 72.7 % with SVMs
- A combination of the three classifiers using logistic regression is also proposed.

Fill in the Blanks for Mid Term Exams

- Q.1 Natural languages are _____ from constructed and formal languages such as those used to program computers or to study logic.
- Q.2 The technique of discovering of word structure is called _____ parsing.
- Q.3 Depending on the means of communication, _____ are spelled out via graphemes or realized through phonemes.
- Q.4 Tokens are _____ words.
- Q.5 In morphology and syntax, a _____ is a morpheme that has syntactic characteristics of a word, but depends phonologically on another word or phrase.
- Q.6 _____ also known as word segmentation, is the fundamental step of morphological analysis and a prerequisite for most language processing applications.
- Q.7 The citation form of a lexeme is also called its _____.

UNIT - II

2

Syntax Analysis

2.1 : Parsing Natural Language

Q.1 Explain the problems by using context free grammar for syntactic analysis of natural language.

▣ Ans. :

- Parsing of a natural language is a process of determining syntactic structure of the text by analyzing the words based on underlying grammar.
- It identifies the information which is not explicitly given in the input sentence.
- For doing this task it requires some additional information i.e. grammar of a language.
- Context Free Grammar (CFG) is used to write the rules of syntax for ex. :

$S \rightarrow NP VP$

$NP \rightarrow \text{'Atul'} \mid \text{'pockets'} \mid DN \mid NP PP$

$VP \rightarrow V NP \mid VP PP$

$V \rightarrow \text{'bought'}$

$D \rightarrow a$

$N \rightarrow \text{'shirt'}$

$PP \rightarrow \text{'PNP'}$

$P \rightarrow \text{'With'}$

Where, V = Verb, NP = Noun phrase, VP = Verb phrase

PP = Prepositional phrases

- Typically the rules are stated in the form $X \rightarrow w$; Where X is a part of speech for word w which is a generated terminal symbol.

For ex. : In rule $V \rightarrow \text{'saw'}$, V is a part of speech or preterminal symbol which generates verb saw.

- The sentence 'Atul bought a shirt with pocket's can have two possible forms based on rules stated in CFG.
- In the first parse form pockets can be considered as a currency to buy shirt and in second form the sentence can be of meaning Atul is purchasing shirts with pockets.

1st parse
 (S (NP Atul)
 (VP (VP (V bought)
 (NP (D a)
 (N shirt)))
 (PP (P with)
 (NP pockets))))))

2nd parse
 (S (NP John)
 (VP (V bought)
 (NP (NP (D a)
 (N shirt))
 (PP (P with)
 (NP pockets))))))

- From the above example we understand that writing all possible syntactic formations in a language is difficult and complex task as we cannot list down a single form. All the possibilities need to be explored based on the part of speech tags mentioned for a particular word.
- It is also difficult to mention lexical properties of a particular word, which is a first knowledge acquisition problem which is a process of extracting structuring and organizing knowledge.
- Second, problem of knowledge acquisition arises from the fact that, it is not only sufficient to know the syntactic rules of a language but it is also required to understand which analysis is most feasible for the sentence in terms of meaning due to ambiguity in syntactic analysis.
- For ex : Consider CFG

N → NN

N → 'natural' | 'language' | 'processing' | 'book'

- If we consider input sentence as natural language processing, two possible ambiguous parses can be generated.

Parse 1
 (N (N (N natural)
 (N language))
 (N processing))

Parse 2
 (N (N natural)
 (N (N language)
 (N processing)))

- Parse 1 leads to the meaning processing of natural language and parse 2 leads to meaning natural way to do language processing.
- This ambiguity problem is a limitation of CFG.

2.2 : Treebanks : A Data Driven Approach of Syntax

Q.2 How to knowledge acquisition problems of context free grammar is addressed by tree bank approach.

OR Write a note on Treebank approach of syntax analysis.

Ans. :

- Due to ambiguous nature of a language two knowledge acquisition problems arise in syntactic analysis
 1. Due to difficulty in exploration of all the possibilities based on part of speech tags mentioned for a particular word.
 2. Due to difficulty in understanding the feasible meaning of sentence due to ambiguity in syntactic analysis.
- To address these problems a data driven approach called as **treebank** can be adapted.
- Treebank is collection of sentences.
- Each sentence in treebank contains complete syntax analysis.

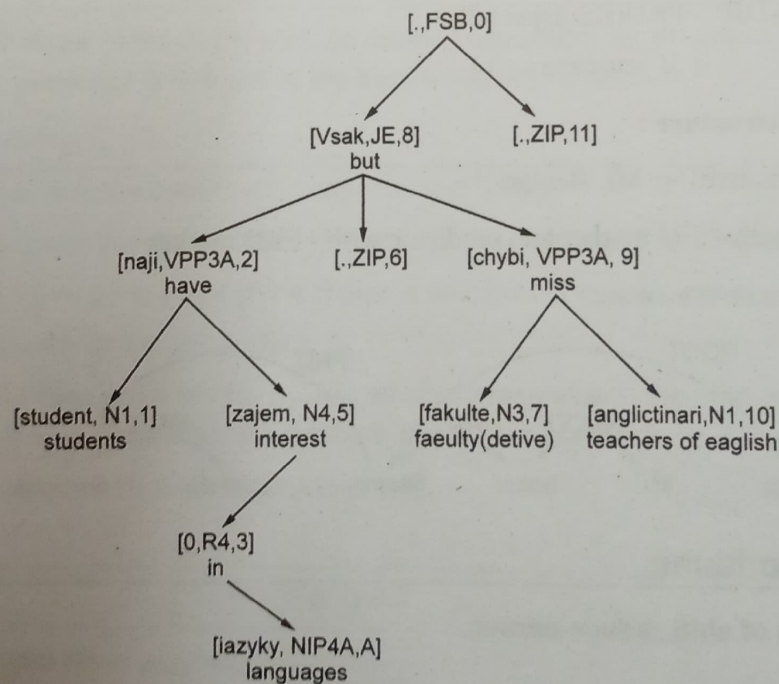
- Human expert provides feasible analysis of the sentence. Annotation guidelines are provided before annotation process.
- Treebank contains annotations of syntactic structure for large set of sentences.
- Supervised machine learning techniques are used to train the parser from the training data extracted from tree banks.
- The first knowledge acquisition problem is addressed by providing syntactic analysis directly instead of grammar.
- The second knowledge acquisition problem is solved as for each sentence in a treebank the most feasible syntactic analysis is provided.
- Using supervised machine learning techniques are used to learn scoring function for all possible syntax analysis.

2.3 : Representation of Syntactic Structure

Q.3 Explain how syntax analysis is done using dependency graphs.

Ans. :

- Dependency graphs connects head of a phrase to its dependents.
- According to definition of Co NLL, in dependency graph nodes are words of input sentence and arcs are binary relations from head to dependent.
- It labelled dependency parsing a label is assigned to each dependency relation between head and dependent word.
- The example dependency graph for Czech sentence from Prague Dependency Treebank is shown in Fig. Q.3.1.



The students are interested in languages, but the faculty is missing teachers of English

Fig. Q.3.1 : An example of a dependency graph syntax analysis for a Czech sentence taken from Prague Dependency Treebank. Each node in the graph is a word, its part of speech, and the not of the word in the sentence, for example [fakulte, N3, 7] is the seventh wood in the sentence with F tag N3, which also tells us that the word has dative case. The node [# , ZSB,0] is the root node of dependency tree. The English equivalent is provided for each node

- It is observed that dependency analysis make minimal assumption about syntactic structure for avoiding annotation of hidden structure like empty elements to represent missing arguments of predicates.

Q.4 Explain the concept of projectivity in dependency analysis.

OR What is projective dependency tree.

▣ **Ans. :** Projectivity is the constraint on syntactic analysis due to effect of linear order of words on dependencies between words.

- The example shown in Fig. Q.4.1 shows the english sentence with the requirements of crossing dependencies.

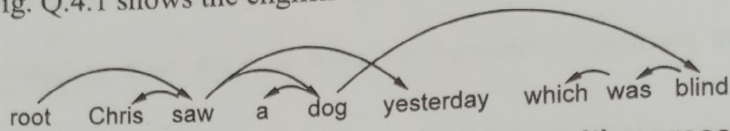


Fig. Q.4.1 : An unlabeled nonprojective dependency tree with a crossing dependency

Q.5 Explain syntax analysis using phrase structure tree.

▣ **Ans. :**

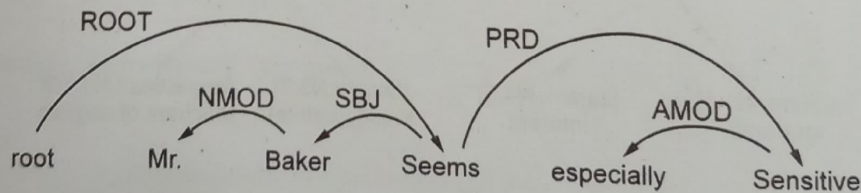
- Phrase structure syntax analysis is based on the concept that the sentences can be divided in constituents and larger constituents are formed by merging smaller ones.
- For ex. : Consider the sentence from penn Treebank "Mr. Baker especially sensitive" The predicate argument structure can be written as,

(S (NP – SBJ (NNP Mr.)
 (NNP Baker))
 (VP (VBZ seems)
 (ADJP – PRD(RB especially)
 (JJ sensitive)))))

Predicate - argument structure :

Seems ((especially (sensitive)) ((Mr. Banker)))

- The subject is marked with - SBJ marker and predicate with - PRD marker.
- The dependency tree can be drawn as,



2.4 : Parsing Algorithms

Q.6 Explain the working of shift reduce parser.

▣ **Ans. :**

- In parsing for a given input string we need to do right most derivation of grammer.
- In shift reduce parsing the concept of pushdown automaton (PDA) is used.
- PDA is an automaton that uses stack.

- The shift reduce parser has two steps
 - Shift step** : In this the input stream is advanced by one symbol. The shifted symbol is considered as single node parse tree.
 - Reduce step** : It applies completed grammar rule to recent parse trees and combine them together as one tree with new root symbol.
- The working of shift reduce parsing algorithm is shown in Fig. Q.6.1.

Parse Tree	Stack	Input	Action
		a and b or c	Init
a	a	and b or c	Shift a
(N a)	N	and b or c	Reduce $N \rightarrow a$
N(a) and	N and	b or c	Shift and
(N a) and b	N and b	or c	Shift b
(N a) and (N b)	N and N	or c	reduce $N \rightarrow b$
(N (N a) and (N b))	N	or c	reduce $N \rightarrow a$
(N (N a) and (N b)) or	N or	c	shift or
(N (N a) and (N b)) or c	N or c		Shift c
(N (N a) and (N b)) or (N c)	N or N		Reduce $N \rightarrow c$
(N (N (N a) and (N b)) or (N c))	N		Reduce $N \rightarrow N$ or N
(N (N (N a) and (N b)) or (N c))	N		Accept!

Fig. Q.6.1 : The individual steps of the shift-reduced parsing algorithm for the input a and b or c for the grammar G defined at the beginning of this section

- The algorithm is defined as Fig. Q.6.2.

<ol style="list-style-type: none"> Start with an empty stack and the buffer containing the input string. Exit with success if the top of the stack contains the start of the grammar and if the buffer is empty. Choose between the following two steps (if the choice is ambiguous, choose one based on an oracle) : <ul style="list-style-type: none"> Shift a symbol from the buffer onto the stack. If the top k symbols of the stack are $\alpha_1, \dots, \alpha_k$, which corresponds to the right hand side of a CFG rule $A \rightarrow \alpha_1, \dots, \alpha_k$, then replace the top k symbols with the left-hand side nonterminal A. Exit with failure if no action can be taken in previous step. Else, go to step 2.
--

Fig. Q.6.2

Q.7 Explain hypergraphs and chart parsing.

OR Explain CYK algorithm.

▣ **Ans. :**

- CFG requires the use of database.
- Due to linear parsing technique in CFG in the worst case ran time of algorithm is exponential in the grammar size.

- To address this instead of left to right parsing statistical parser which search the space for possible sub trees is used.
- As shown in below example the example grammar can be rewritten in which right hand side contain only two non-terminals.

Example G
 $N \rightarrow N \text{ 'and' } N$
 $N \rightarrow N \text{ 'or' } N$
 $N \rightarrow \text{ 'a' } | \text{ 'b' } | \text{ 'c' } |$

New Gc
 $N \rightarrow N N^{\wedge}$
 $N^{\wedge} \rightarrow \text{ 'and' } N$
 $N \rightarrow N N_{\vee}$
 $N_{\vee} \rightarrow N \text{ 'or' } N$
 $N \rightarrow \text{ 'a' } | \text{ 'b' } | \text{ 'c' } |$

- This can be further made compact by linking input sentence into spans 0 a 1 and 2 b 3 or 4 c 5 etc. i.e. string a is in span 0, 1.
 b or c is in span 2, 5.
- So a new Grammar G_f is written using this concept as below :

New Gf
 $N [0, 5] \rightarrow N [0, 1] N^{\wedge} [1, 5]$
 $N [0, 3] \rightarrow N [0, 1] N^{\wedge} [1, 3]$
 $N^{\wedge} [1, 3] \rightarrow \text{ 'and' } [1, 2] N [2, 3]$
 $N^{\wedge} [1, 5] \rightarrow \text{ 'and' } [1, 2] N [2, 5]$
 $N [0, 5] \rightarrow N [0, 3] N_{\vee} [3, 5]$
 $N [2, 5] \rightarrow N [2, 3] N_{\vee} [3, 5]$
 $N_{\vee} [3, 5] \rightarrow \text{ 'or' } [3, 4] N [4, 5]$
 $N [0, 1] \rightarrow \text{ 'a' } [0, 1]$
 $N [2, 3] \rightarrow \text{ 'b' } [2, 3]$
 $N [4, 5] \rightarrow \text{ 'c' } [4, 5]$

- The parse tree that starts from a start nonterminal and spans complete string is shown in Fig. Q.7.1.

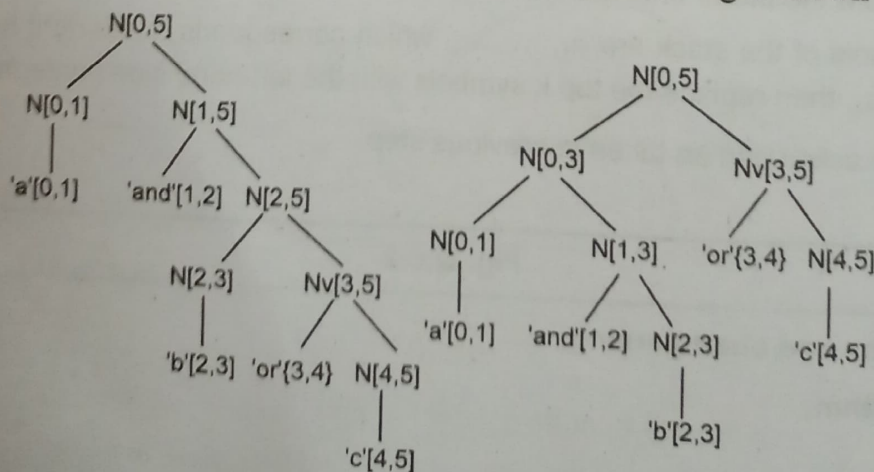


Fig. Q.7.1 : Parse trees embedded in the specialized CFG for a particular input string. The nodes with the same label, such as $N[0,5]$, $N[0,1]$ and $[1,2]$, $N[2,3]$ and $N_{\vee}[3,5]$ can be merged to form a hypergraph representation of all parses for the input

- To construct this specialized CFG following are steps :
- First rules generating lexical items are considered

For ex. $N [0, 1] \rightarrow 'a' [0, 1]$

$N [2, 3] \rightarrow 'b' [2, 3]$

$N [4, 5] \rightarrow 'c' [4, 5]$

- The pseudocode is written as shown in Fig. Q.7.2.

```

for i = 0 ..... n do
    if  $N \rightarrow x$  with score  $s$  for any  $x$  spanning  $i, i + 1$  exists then
        add specialized rule  $N [i, i + 1] \rightarrow x [i, i + 1]$  with score  $s$ 
        written as  $N[i, i + 1] : s$ 
    end if
end for

```

Fig. Q.7.2

- In the next step specialized rules are created recursively. For ex. If $Y[i, k]$ and $Z[k, j]$ are left hand sides of previously created rules then rule $X \rightarrow YZ$ can be converted to

$X [i, j] \rightarrow Y [i, k] Z [k, j]$

- Score S is assigned to each non-terminal span.
- The algorithm is shown in Fig. Q.7.3 and is known as CKY (Cocke, Kasami and younger) algorithm.

```

for j = 2 .... n do
    for i = j - 1 ... 0 do
        for k = i + 1 .... j do
            if  $Y [i, k] : s_1$  and  $Z[k, j] : s_2$  are in the specialized grammar then
                if  $X \rightarrow YZ$  with score  $s$  exists in the original grammar then
                    add specialized rule  $X[i, j] \rightarrow Y[i, k] Z[k, j]$  with score  $s + s_1 + s_2$ 
                    keep only the highest scoring rule :  $X[i, j] \rightarrow \alpha$ 
                end if
            end if
        end for
    end for
end for
end for

```

Fig. Q.7.3

Q.8 Explain working of minimum spanning trees and dependency parsing.

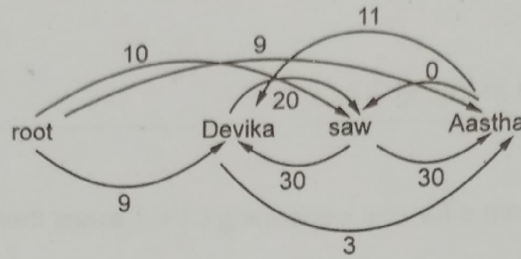
▣ **Ans. :**

- The Minimum Spanning Tree (MST) corresponds to the optimum branching problem in directed graphs. Which are rooted and does not have cycles.
- The basic prerequisite is all the dependency links between the words must have score.

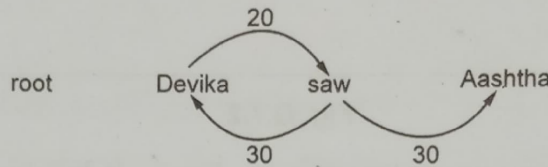
- To understand the working of MST consider the input sentence as

Devika saw Aastha

- The fully connected graph of this sentence can be drawn as follows,



- The scoring function is used to assign the weights to the edges.
- The algorithm starts with finding the incoming edge with highest score.



- We have a cycle in this graph. We can combine the cycle into single node and recalculate edge weight.
- The edge weight from each node to this newly combined node is computed.
- Also record the maximum weight of a node in this combination.
- For example the incoming edge.

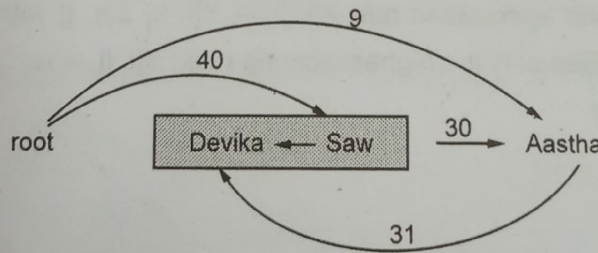
root → **saw → Devika** is having weight $10 + 30 = 40$

Whereas root → **Devita → saw** is having weight $9 + 20 = 29$

Aastha → **saw → Devika** is having weight $0 + 30 = 30$

and Aastha → **Devika → saw** is having weight $11 + 20 = 31$

So the graph can be viewed as

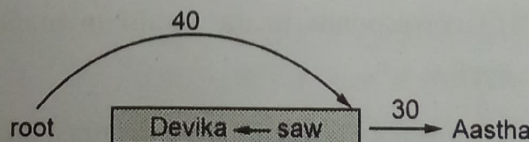


- The MST algorithm is recursively applied to this graph and best incoming edges to each word are found.
- So in the next iteration

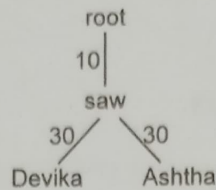
root → Aastha → **Devika → saw** will have weight $9 + 31 = 40$

and root → **Devika → saw** → Aastha will have weight $40 + 30 = 70$

- So the best nodes chosen are shown in below graph



And we get a highest scoring dependency parse as



2.5 : Models for Ambiguity Resolution in Parsing

Q.9 How ambiguity resolution in parsing is done with the help of probabilistic context free grammars.

Ans. :

- Consider the example sentence “Atul bought a shirt with pockets” explain a in θ .
- Please refer to the CFG and parse trees of this sentence from θ_1 .
- To resolve the ambiguity of such type, one way is to assign the probabilities to the rules in CFG.
- Due to this the CFG is known as Probabilistic Context Free Grammar or PCFG. For example for the rule $N \rightarrow \alpha$ the probability can be defined as $P(N \rightarrow \alpha/N)$ such that rule probability is stated at left hand side.
- Due to this assignment when non-terminal is expanded, probability distribution is done among all expansions of non-terminals. i.e.

$$1 = \sum_{\alpha} P(N \rightarrow \alpha)$$

- So for the example sentence the probability distribution can be viewed as

$$S \rightarrow NP VP (1.0)$$

$$NP \rightarrow \text{'Atul'} (0.1) \mid \text{'pockets'} (0.1) \mid DN (0.3) \mid NP PP (0.5)$$

$$VP \rightarrow V NP (0.9) \mid VP PP (0.1)$$

$$V \rightarrow \text{'bought'} (1.0)$$

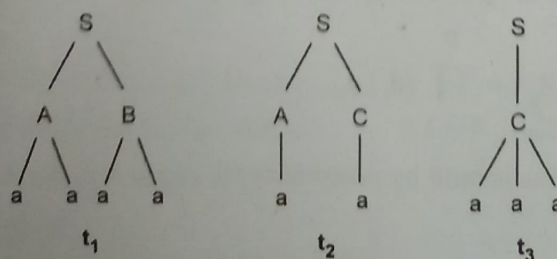
$$D \rightarrow \text{'a'} (1.0)$$

$$N \rightarrow \text{'shirt'} (1.0)$$

$$PP \rightarrow P NP (1.0)$$

$$P \rightarrow \text{'with'} (1.0)$$

- From the assigned probabilities we can observe that the decision should be taken from two rules $NP \rightarrow NP PP$ and $VP \rightarrow VP PP$ and as probability of $NP \rightarrow NP PP$ is having higher probability it generates a more feasible sentence.
- Consider the following example of derivation of rule probabilities from treebank. The figure below shows treebank with three trees t_1 , t_2 and t_3 .

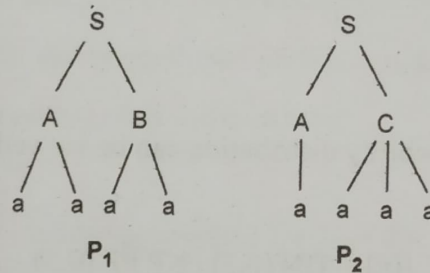


Assume that t_1 occur 10 times, t_2 : 20 times and t_3 : 50 times in treebank.

PCFG for this can be given as

$\frac{10}{10 + 20 + 50} = 0.125$	$S \rightarrow A B$
$\frac{20}{10 + 20 + 50} = 0.25$	$S \rightarrow A C$
$\frac{50}{10 + 20 + 50} = 0.625$	$S \rightarrow C$
$\frac{10}{10 + 20} = 0.334$	$A \rightarrow a a$
$\frac{20}{10 + 20} = 0.667$	$A \rightarrow a$
$\frac{20}{20 + 50} = 0.285$	$B \rightarrow a a$
$\frac{50}{20 + 50} = 0.714$	$C \rightarrow a a a$

- If we consider input a a a a, we can generate two parses.



Where, $P_1 = 0.125 * 0.334 * 0.285 = 0.01189$
 $P_2 = 0.25 * 0.667 * 0.714 = 0.119$

So P_2 is most feasible for parsing.

Q.10 Explain generative models for parsing for ambiguity resolution.

▣ Ans. :

- A parse tree is typically built by sequence of decisions.
- Based on the CFG rules there can be multiple derivations.
- Consider each such derivation as, $D = d_1, \dots, d_n$.
- Parser has to choose from these derivations, the most feasible one.
- Let us consider input sentence as x and output parse tree need to be generated as y .
- For each derivation of parse tree the probability can be assigned as

$$P(x, y) = P(d_1, \dots, d_n) = \prod_{i=1}^n (d_i | d_1, \dots, d_{i-1})$$

- In this equation a partial parse tree is built by probability $(d_i | d_1, \dots, d_{i-1})$ which is known as history.

- These histories are grouped into conditional classes by function ϕ as :

$$P(d_1, \dots, d_n) = \prod_{i=1}^n (d_i | \phi(d_1, \dots, d_{i-1}))$$

- History $H_i = d_1, \dots, d_{i-1}$ for all x, y is represented by finite set of feature functions k

$$\phi_1(H_i), \dots, \phi_k(H_i)$$

$$P(d_1, \dots, d_n) = \prod_{i=1}^n [d_i | \phi_1(H_i), \dots, \phi_k(H_i)]$$

Q.11 Explain global linear discriminative model ambiguity resolution parsing.

OR Explain global linear model.

▣ **Ans. :**

- Discriminative model is developed by Collins for creating simple framework for describing discriminative approaches, which is also called as global linear model.
- Consider X as set of inputs and Y as output which is sequence of POS tags as parse tree.
- Each input $x \in X$ and $y \in Y$ is mapped to d -dimensional feature vector $\phi(x, y)$. Each dimension is a real number.

- Weight is assigned to each feature is $\phi(x, y)$ by weight parameter vector $w \in \mathbb{R}^d$.

$$\phi(x, y) \cdot w = \text{Score of } (x, y)$$

If the score is higher y is the most feasible output for x .

- Possible outputs y from x are computed from function $\text{GEN}(x)$.
- The highest scoring candidate y^* from $\text{GEN}(x)$ is computed as

$$F(x) = \underset{y \in \text{GEN}(x)}{\text{argmax}} P(y | x, w)$$

- Conditional random field C (RF) compute conditional probability as

$$\log p(y | x, w) = \phi(x, y) \cdot w - \log \sum_{y' \in \text{GEN}(x)} \exp[\phi(x, y') \cdot w]$$

- Global linear model is stated as

$$F(x) = \underset{y \in \text{GEN}(x)}{\text{argmax}} \phi(x, y) \cdot w$$

Q.12 Explain the original perceptron learning algorithm for ambiguity resolution in parsing.

▣ **Ans. :**

- Perceptron is a single layered neural network.
- It process a example at a time.
- The weight adjustment is done of weight parameter vector.
- This vector is applied to input to generate the related output.
- The features present in the truth are "recognized."
- Consider a training set with in examples.

- As shown in Fig. Q.12.1 Algorithm 12.1 the original perceptron learning Q.12.1 Algorithm works like below.

Algorithm 12.1 : The original perceptron learning algorithm
Inputs : Training data $\{ (x_1, y_1), \dots, (x_m, y_m) \}$; number of iterations T
Initialization : Set $w = 0$
Algorithm :

1. **for** $t = 1, \dots, T$ **do**
2. **for** $i = 1, \dots, m$ **do**
3. Calculate y_i , where $y_i' = \operatorname{argmax}_{y \in \operatorname{GEN}(x)} \Phi(x_i, y) \cdot w$
4. **if** $y_i' \neq y_i$ **then**
5. $w = w + \Phi(x_i, y_i) - \Phi(x_i, y_i')$
6. **end if**
7. **end for**
8. **end for output** : The updated weight parameter vector w

Fig. Q.12.1 : Algorithm 12.1

1. Weight parameter w is initialized to 0.
 2. Iteration is carried out on m training examples.
 3. Set of candidates $\operatorname{GEN}(x)$ is generated for each x .
 4. The most feasible candidate i.e. the candidate with maximum score according to w is selected.
 5. w is updated by increasing weight values of features in truth and decreasing weight value for features appearing in top candidate.
- The problem faced by this algorithm is of overfitting during incremental weight update due to which unseen data is not classified properly.
 - The algorithm is not suitable for linearly inseparable training data.

Q.13 Explain voted perceptron algorithm for ambiguity resolution in parsing.

▣ **Ans. :**

- As shown in Fig. Q.13.1 Algorithm 13.1 the voted perceptron algorithm, proposed by Freund and Schapire and works as follows :
 - Instead of a single weight vector W , the learning process considers all intermediate weight vectors.
 - In classification phase-these intermediate vectors are used to vote for the answer.
 - Good prediction vector generally survive for long time and have larger weight in the vote.
 - In training phase count C_i counts the survival of weight parameter vector (w_i, C_i) in training.
 - If top candidate is not in truth C_{i+1} is initialized 1, to generate an updated weight vector (w_{i+1}, C_{i+1}) .
 - While this original C_i and weight vector (w_i, C_i) are stored.
 - This algorithm is more stable - than original perceptron.

Algorithm 13.1 : The voted perception algorithm

Training phase

Input : Training data $\{ (x_1, y_1), \dots, (x_m, y_m) \}$; number of iterations T

Initialization : $k = 0, w_0 = 0, c_1 = 0$

Algorithm :

for $t = 1, \dots, T$ do

 for $i = 1, \dots, m$ do

 Calculate y'_i , where $y'_i = \operatorname{argmax}_{y \in \text{GEN}(x)} \Phi(x_i, y) \cdot w_k$

 if $y'_i \neq y_i$ then

$c_k = c_k + 1$

 else

$w_{k+1} = w_k + \Phi(x_i, y_i) - \Phi(x_i, y'_i)$

$c_{k+1} = 1$

$k = k + 1$

 end if

 end for

end for

output : A list of weight vectors $((w_1, c_1), \dots, (w_k, c_k))$

Prediction Phase

Input : The list of weight vectors $((w_1, c_1), \dots, (w_k, c_k))$, an unsegmented sentence x .

Calculate :

$$y^* = \operatorname{argmax}_{y \in \text{GEN}(x)} \left(\sum_{i=1}^k c_i \Phi(x, y) \cdot w_i \right)$$

Output : The voted top ranked candidate y^* .

Fig. Q.13.1 : Algorithm 13.1

Q.14 Explain averaged perceptron algorithm for ambiguity resolution in parsing.

▣ **Ans. :**

- The averaged perceptron algorithm reduces space and time complexities.
- As shown in Fig. Q.14.1 Algorithm 14.1 instead of w -the averaged weight parameter vector γ on m training examples is used.
- γ can be defined as

$$\gamma = \frac{1}{mT} \sum_{i=1, \dots, m, t=1, \dots, T} w^{i,t}$$

- Accumulating parameter vector σ is maintained and updated using w .

Algorithm 14.1 : The averaged perceptron learning algorithm

Input : Training Data $\langle (x_1, y_1), \dots, (x_m, y_m) \rangle$; number of iterations T

Initialization : Set $w = 0, \gamma = 0, \sigma = 0$

Algorithm :

```

for t = 1, ..... , T do
    for i = 1, ..... , m do
        Calculate  $y_i$ , where  $y_i = \operatorname{argmax}_{y \in \text{GEN}(x)} \Phi(x_i, y) \cdot w$ 
    end for
    if  $y_i \neq y_i$  then
         $w = w + \Phi(x_i, y_i) - \Phi(x_i, y_i)$ 
    end if
     $\sigma = \sigma + w$ 
end for
end for
output : The averaged weight parameter vector  $\gamma = \sigma / (mT)$ 
words, are used as heads and dependent

```

Fig. Q.14.1 : Algorithm 14.1

2.6 : Multilingual Issues

Q.15 What are the different issues faced by a parser in terms of tokenization, case and encoding ?

▣ **Ans. :**

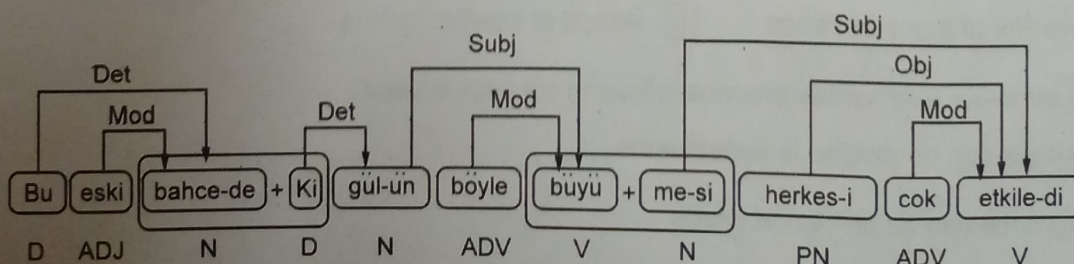
- Typically in a grammar a word has a specific definition.
- But the definition of a word may be different for different parsers or treebanks. For ex. In English language two tokens are separated by space, but in parser for english word today's or there's are considered as two independent tokens i.e. today and 's, there and 's.
- Second issue is with uppercase and lowercase words.
- If we simply convert all the uppercase of treebank to the lowercase we may loose some useful information. For ex : In proper nouns like India the first letter will always in a uppercase form but if we lowercase all the letters then in the training data the proper nouns may look like a simple verbs.
- Solution of this problem can be selective lower casing. Also low count token can be replaced by a pattern. For ex. : India appears twice it is replace of by in dates same is applicable to URL's or IP addresses, etc.
- The next issue is with the encoding style of the 'language script'.
- If the script is not in ASCII encoding then the encoding needs to be handled.
- Care should be taken-that data used by parser is converted to the encoding of tree bank and vice versa.
- The text data may take different encoding styles (based on the language script) such as UTF-8, GB, BIG5, for Chinese text.

Q.16 Explain the role of word segmentation in parsing.**Ans. :**

- In the languages like Chinese the words are not separated by the identification marks.
The written form of many languages, including Chinese, lack marks identifying words. Given the Chinese text, **北京大学生比赛**, a plausible segmentation would be **北京** (Beijing) **大学生** (university students) / **比赛** (competition) 'competition among university students in Beijing'. However, if **北京大** is taken to mean Beijing University, the segmentation.
- From the above example we can understand the importance of word segmentation.
- It is process by which the character sequence is divided into blocks so that the generated output is meaningful but consists of separate tokens.
- After identification of all the words in a sentence POS tags are assigned to each word and syntax tree for a sentence can be built.
- This task is easy for the languages like English as all the tokens (words) are separated by spaces, but quite challenging for the languages like Chinese, Japanese, etc.
- Various researchers have addressed this problem by different approaches.
- In one approach for Chinese text parsing, the parser assigns word boundaries where non terminals of tree specify word boundaries.
- But the immediate context proves to be more useful in detecting word boundaries.
- In the approach stated by Bar-Hillel, Perles and Shamir, the parser consider states in automata as indices. The input is considered as a directed acyclic graph with a single start point and edges labeled with word and height called as word lattice. This word lattice represent multiple segmentation probabilities of the segmented Chinese text. The most feasible ranked segmentation is chosen by the parser to generate most accurate parse.

Q.17 Explain the role of morphology in syntactic analysis.**OR Explain how morphology affects parsing of the languages other than English.****Ans. :**

- Basic components of a word are called morphemes.
- Meaning of the word is derived from combination of meanings of morphemes.
- A word is considered to be a combination of different morphemes contributing to the meaning.
- For ex : Consider a Turkish sentence shown in Fig. Q.17.1 '+' symbol shows the morphemes present in a word. So it is feasible to use morphemes as heads and dependents instead of words.

**Fig. Q.17.1**

- In the languages like Czech and Russian morphemes are also used to show grammatical case, genders, etc.
- For Ex. : In Czech language maximum adjectives form all four genders, seven case markers, all three degrees of comparison and are positive and negative in polarity, resulting in 336 inflected words only for adjectives.
- One solution to resolve the ambiguity in morphology is to assign POS tag for encoding various morphemes. For ex. : POS V...M-3 indicates that each word contain morphemes in 10 dimensions and in this example stem is verb with masculine gender in third person.
- This is accomplished by training sub classifiers for each component of POS tag.

Fill in the Blanks for Mid Term Exam

- Q.1 _____ of a natural language is a process of determining syntactic structure of the text by analyzing the words based on underlying grammar.
- Q.2 _____ contains annotations of syntactic structure for large set of sentences.
- Q.3 _____ analysis is based on the concept that the sentences can be divided in constituents and larger constituents are formed by merging smaller ones.
- Q.4 In shift reduce parsing the concept of _____ is used.
- Q.5 The _____ prerequisite is all the dependency links between the words must have score.
- Q.6 To resolve the ambiguity of in the sentence, one way is to assign the _____ to the rules in CFG.
- Q.7 _____ developed by Collins for creating simple framework for describing discriminative approaches, which is also called as global linear model.
- Q.8 Basic components of a word are called _____.
- Q.9 _____ of the word is derived from combination of meanings of morphemes.

Multiple Choice Questions for Mid Term Exam

- Q.1 Parsing of Natural Language Processing is determination of _____.
- a Semantic structure b Syntactic structure
- Q.2 What is a tree bank ?
- a Collection of words b Collection of sentences
- c Collection of letters d Collection of tokens
- Q.3 What are nodes of dependency graph ?
- a Alphabets of input sentence b Words of input sentence
- Q.4 In phrase structure tree syntax analysis following concept is used :
- a Sentences can be divided in constituents
- b Paragraphs can be divided in words



3.1 : Introduction

Q.1 What is semantic parsing ?

▣ Ans. :

- Semantic parsing is an important phase of Natural Language Processing (NLP) .
- Any document is comprised of set of sentences.
- These sentences are formed by the words arranged according to the rule based grammar as per the constructs in a particular language.
- To extract the meaning out of this syntactically arranged sentences is the objective of semantic parsing.
- Parsing is analysing a text into logical syntactic components and semantics is study of meaning.
- The information pieces are identified and related in semantic parsing.
- The meaningful parts are identified in the text and they are transformed into suitable data structures for higher level task accomplishment.

Q.2 What are types of semantic parsing ?

OR Differentiate between deep semantic parsing and shallow semantic parsing.

▣ Ans. :

- Semantic parsing is about identifying the meaning.
- There are two approaches which can be adapted extracting the meaning out of text.
- Consider domain specific applications like travel reservation, gaming simulation etc.
- In these kind of applications precise and each meaning extraction can be done based on the restricted limits of the particular domain.
- In this approach based on the query or the purpose, output is generated from the meaning representation. This approach is known as **deep semantic parsing**.
- The second approach is not domain specific and is more generic.
- In this approach the task of meaning representation is divided in small pieces. It is known as shallow semantic parsing.
- These pieces are responsible for capturing small manageable components that represent meaning.

- The related sets of meaning representations are created in this case for example extraction of word sense disambiguation followed by predicate argument structure.
- In case of deep semantic parsing there is little or no scope of reusability as every domain is unique and requires different understanding and representation of the concepts.
- In case of shallow semantic parsing as generic meaningful pieces need to be created, it is very difficult to have general purpose ontology and symbols which are shallow from learning point of view but have high reusability irrespective of the applications.

3.2 : Semantic Interpretation

Q.3 What is semantic Interpretation ?

▣ Ans. :

- Semantic interpretation facilitates the joining of different components which define meaning representation of the text.
- This representation, when fed to a computer can be further processed for computational manipulations and search for any application.
- Semantic parsing is the part of semantic interpretation.

Q.4 Explain semantic theory.

▣ Ans. :

- The semantic theory was proposed by katz and fodor in 1963.
- It addresses the fact of meaning representation which is lacking in Chomsky's transformational grammar.
- Following points are stated in semantic theory :
 1. Ambiguous sentences should be taken care of and explained properly. Consider the word 'bank' in the sentence 'I went to bank'. In this sentence ambiguity exist as there can be two meaning of the word bank first is money bank and second is river bank.
 2. Ambiguity resolution of the words depending on context should be done. For example : If the sentence is extended as 'I went to bank to withdraw money' then the theory should be able to disambiguate word 'bank' to extract correct meaning.
 3. The theory should be able to identify meaningless but syntactically correct sentences. For example : Consider the sentence 'Colorless green ideas sleep furiously'. This sentence does not carry any meaning but is syntactically correct.
 4. Many times same semantic content can be represented by different unrelated syntactic structures such sentences should be handled.

Q.5 List and explain all the components in semantic interpretation.

▣ Ans. : The process of semantic interpretation consists of different components to represent a text which can be fed into and processed by computer for undergoing various functionalities like search which are basis of language understanding system.

- Following are the major components in this process

□ 1. Structural ambiguity

- Any sentence is represented by its syntactic structure.
- Syntax and semantics are closely related to each other and we consider that the semantic interpretation is based on underlying syntactic structure.
- So syntactic processing is the first step in semantic interpretation.

□ 2. Word sense

- Many times in any language a same word is used in different meaning depending on the world knowledge.
- For example a word "bank" represents a money bank or it can also possess a meaning river bank.
- Due to inherent intelligence and language vocabulary by humans it is not a difficult for them to understand the meaning of the word expected by speaker or author.
- Consider the examples :
 1. Go to bank to withdraw money.
 2. Fetch some water from river bank.
- It's a easy task for humans to disambiguate the meaning of word bank in above sentences.
- But for a machine this word sense disambiguation is a challenging task and it plays important role in semantic interpretation.

□ 3. Entity and entity resolution

- Any text consists of various entities falling in different categories like person name, locations, quantities, etc.
- Identification of these entities is the major task in semantic interpretation system.
- Named Entity Recognition (NER) is a subtask of information extraction for locating and classifying named entities in the above mentioned categories.
- For example : [Devika]_{Person} bought 100 shares of [Infosys]_{organization} in [2019]_{time}
- Another important task is of co-reference resolution.
- Co-reference occurs when two or more expressions in text refer to same person.
For example : Pranjali said she will sing.
In this sentence proper noun Pranjali and she refers to same person.
- This task is also under information extraction and a major component in semantic interpretation.

4. Predicate argument structure

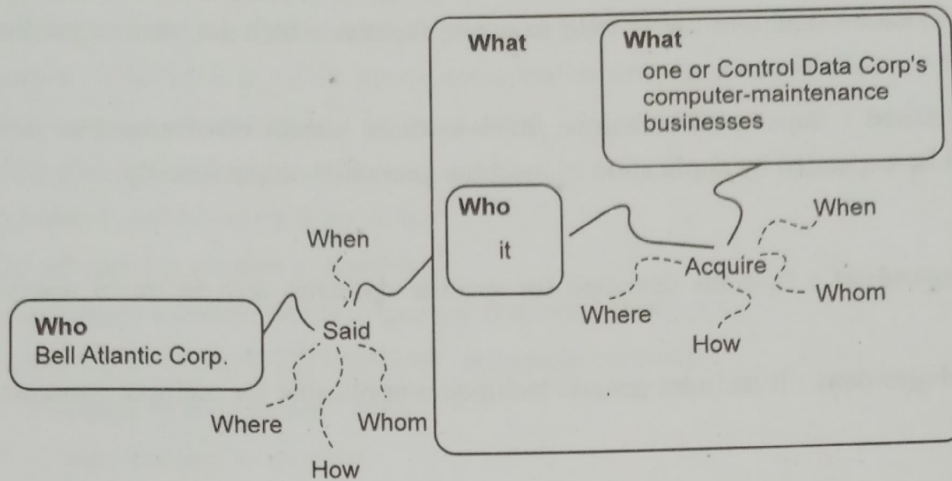


Fig. Q.5.1 : A representation of *who did what to whom, when, where, why and how*

- After finishing all the above mentioned tasks of word sense disambiguation, NER and coreference resolution, next task is to identify the participants in the events.
- It is important to identify what is a role of entity in a particular event.
- This is known as resolving the argument structure of predicate in a sentence.
- It is typically identification of who did what to whom, when, where and how as shown in Fig. Q.5.1.

5. Meaning representation

- This is the last step in semantic interpretation and also called as deep representation.
- It involves building meaning representation which can be used by algorithms for various applications.
- Research is still going on in the area of general purpose representations.
- Much study has been done in domain specific applications.
- For example : Consider a Geoquery domain, the example query can be represented as which river is the longest ?

Answer (x_1 , longest (x , river (x_1)))

3.3 : System Paradigms

Q.6 What are the different approaches for practical simple mentation of semantic interpretation ?

Ans. :

- Based on the diversity of languages and levels of granularity and generality there can be different approaches to handle semantic representation.
- The main constraint is posed by lack of data availability of the languages.
- Due to these constraints some successful approaches for practical implementation of semantic interpretation fall in three categories :

1. System architectures

1. **Knowledge based** : These systems are designed using predefined set of rules for finding the solution to a particular problem.
2. **Unsupervised** : Solution to a particular application is found by minimum human involvement.

3. **Supervised** : This technique involves model training by application of different machine learning algorithms. Feature functions are created to create features which are used to predict labels to handle unseen data.
4. **Semi supervised** : Supervised technique involves more human involvement so is expensive. So the data set can be expanded by application of machine generated output directly.

□ 2. Scope

1. **Domain dependent** : Systems designed for specific domains such as travel reservations or football coaching.
2. **Domain independent** : It includes generic techniques applicable for multiple domains.

□ 3. Coverage

1. **Shallow** : In this approach intermediate representation is generated to be used by machine.
2. **Deep** : Through this approach representation which is created, directly used by machine.

Q.7 Write a note on availability of the resources for word sense disambiguation.

OR What are the different resources for word sense disambiguation in language understanding ?

□ **Ans. :**

- In case of natural language understanding the availability of the corpus is the crucial factor.
- Uptil now the inadequacy of the tagged sense data is the major issue faced by language understanding systems.
- In case of word sense disambiguation the resources are evolved from machine readable dictionaries to the use of word net in the recent systems.
- In early days Longman Dictionary of Contemporary English (LDOCE) and Roget's thesaurus were the resources used for task of disambiguation.
- In late 1980's WordNet was developed which contributed significantly in this task.
- WordNet is a powerful resource which contains lexical database of word senses with multiple parts of speech of a language.
- It also handles different relationships like hypernymy, homonymy, metonymy etc. connecting different words.

3.4 : Word Sense

Q.8 Explain how word sense disambiguation can be done by rule based systems using Lesk algorithm.

□ **Ans. :**

- In early days word sense disambiguation used to done using dictionary sense definitions.
- These are the man made resources and now this has become historical data which is available in archived publication.
- Despite of the fact that this information is not available today, some algorithms and techniques are still useful.
- The first and very simple algorithm was proposed by lesk and these first generation algorithms are based on computerized dictionaries.
- Lesk algorithm is further simplified as shown in Fig. Q.8.1 Algorithm 3.1.

Algorithm 3.1 : Pseudocode of the simplified Lesk algorithm

The function COMPUTEOVERLAP returns the number of words common to the two sets

Procedure : SIMPLIFIED_LESK (word, sentence) **returns** best sense of word

1. *best-sense* ← most frequent sense of word
2. *max-overlap* ← 0
3. *context* ← set of words in sentence
4. **for all** *sense* ∈ senses of word **do**
5. *signature* ← set of words in gloss and examples of *sense*
6. *overlap* ← COMPUTEOVERLAP (*signature*, *context*)
7. **if** *overlap* gt *max-overlap* **then**
8. *max-overlap* ← *overlap*
9. *best-sense* ← *sense*
10. **end if**
11. **end for**
12. **return** *best-sense*

Fig. Q.8.1 : Algorithm 3.1

Q.9 Explain the working of dictionary based algorithm by Yarovsky using Roget's Thesaurus categories.

OR Explain the algorithm for disambiguating words in Roget's thesaurus categories.

Ans. :

- This dictionary based algorithm uses Groliers Encyclopedia. It classifies unseen words in 1042 categories based on statistical analysis of 100 word list in alphabetical order.
- The algorithm is as shown in Fig. Q.9.1.

1. Collect contexts for each of the *Roget's Thesaurus* categories.
2. Determine weights for each of the salient words in the context.

$$\frac{P(w_i | RCat)}{P(w_i)}$$

3. Use the weights for predicting the appropriate category of the word in the test corpus.

$$\arg \max_{RCat} \sum_w \log \frac{P(w_i | RCat) P(RCat)}{P(w_i)}$$

Fig. Q.9.1 : Algorithm for disambiguating words into *Roget's Thesaurus* categories

- The algorithm works as follows :
 - In step 1 contexts are collected.
 - In the next step weight for each word is computed. For this 50 words are used for each side of context word.

$P(w / \text{RCat}) = \text{Probability of occurrence of word } w \text{ in context of Roget's Thesaurus category RCat.}$

- In the final step unseen words in test set are classified into category which maximum weight.

Q.10 Explain Structural Semantic Interconnections (SSI) algorithm for word sense disambiguation.

▣ **Ans. :**

- The algorithm is designed by Navigli and Velardi and uses graphical representation for word sense disambiguation.
- Consider the example shown in Fig. Q.10.1 for the noun bus with a vertical sense and connector sense.
- The algorithm fetches the information from following resources to understand the context of given word
 - WordNet
 - Domain labels
 - Annotated corpora to generate semantic graph.
- The SSI algorithm consists of two steps :
 1. An initialization step
 2. An interactive step.

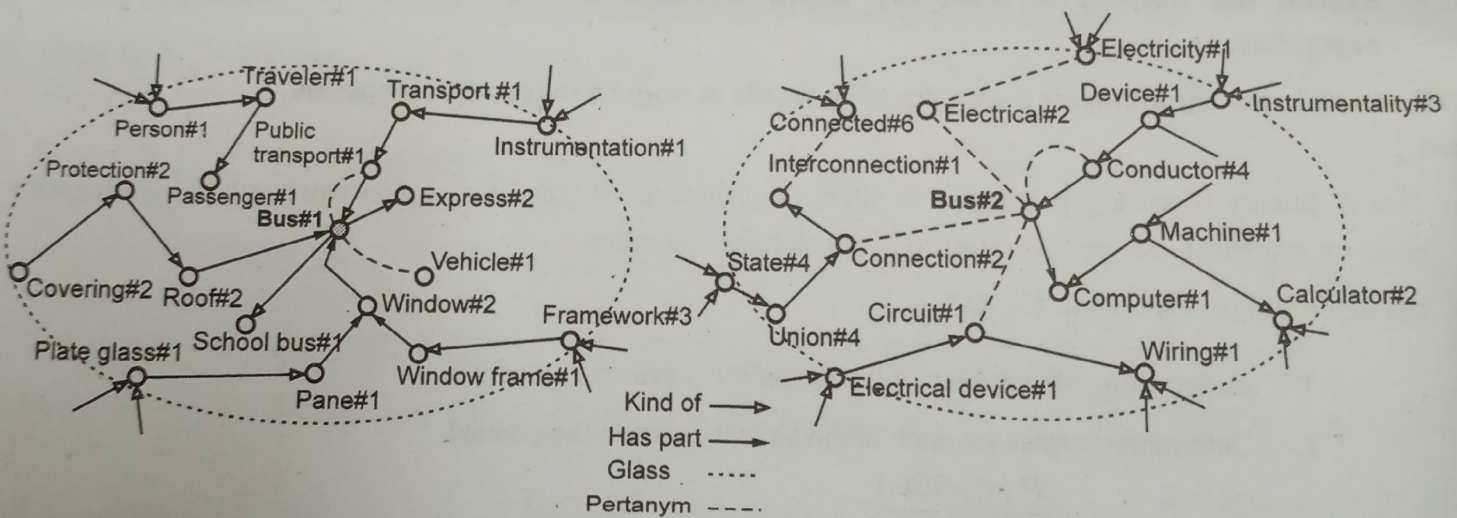


Fig. Q.10.1 : The graphs for sense 1 and 2 of the noun bus as generated by the SSI algorithm

- The algorithm disambiguate the word of a particular context, till further disambiguation is not possible.
- The algorithm has notations as shown in Fig. Q.10.2.

Notation :

- T Notation (the *lexical context*) is the list of terms in the context of the term t to be disambiguated. $T = [t_1, t_2, \dots, t_n]$.
- $S_1^t, S_2^t, \dots, S_n^t$ are structural specifications of the possible concepts (or, senses) of t.
- I (the *semantic context*) is the list of structural specifications of the concepts associated with each of the terms in $T \setminus \{t\}$ (except t). $I = [S^{t_1}, S^{t_2}, \dots, S^{t_n}]$, that is, the *semantic interpretation of T*.

- G is the grammar defining the various relations between the structural specifications (or **semantic interconnections**) among the graphs.
- Determine how well the structural specifications in I match that $S_1^t, S_2^t, \dots, S_n^t$ using G.
- Select the best matching S_1^t .

Fig. Q.10.2

- It works as follows
 - Great set of pending terms in context
 $P = \{t : | S^t = \text{null}\}$
 - Use I to disambiguate P in each iteration.
 - As the algorithm proceeds, which each iteration, one term in P is disambiguated and removed from pending list.
 - It stops when no more terms remain for disambiguation.
 - Output I is updated with sense of t.
- Select terms t in P having semantic interconnections which one sense of S of t and one or more senses in I.
- Function f I (S, t) is computed to determine likelihood of S to be correct interpretation of t

$$f_1(S, t) = \begin{cases} P(\{\Psi(S, S') \mid S' \in I\}) & \text{if } S \in \text{Senses}(t) \\ 0 & \text{Otherwise} \end{cases}$$

Where,

Senses (t) = Senses associated with t

$$\Psi(S, S') = P'(\{\omega(e_1 \cdot e_2 \dots e_n) \mid S \xrightarrow{e_1} S_1 \xrightarrow{e_2} \dots \xrightarrow{e_{n-1}} S_{n-1} \xrightarrow{e_n} S'\})$$

Function P' of weights ω connecting S and S' where S and S' are semantic graphs are e to e_n are edges.

- At last CFG
 - G = (E, N, SG, PG) includes meaningful semantic patterns
 - where E = { $e_{\text{kind-of}}, e_{\text{has-kind}}, e_{\text{part-of}}, \dots$ } are the edge labels.
 - N = { $S_{\text{or}}, S_s, S_g, S_1, S_2, \dots, E_1, E_2, \dots$ } are the nonterminal symbols.
 - S_G = Start symbol of graph.
 - $P_G = \{S_G \rightarrow S_s \mid S_g, S_s \rightarrow S_1 \mid S_2 \mid S_3, S_1 \rightarrow E, S, \mid E_1, E_1 \rightarrow e_{\text{kind-of}} \mid e_{\text{part-of}}, S_g \rightarrow e_{\text{glass}} S_5 \mid S_4 \mid S_5, \dots\}$

Q.11 Explain how Word Sense Disambiguation (WSD) can be done by supervised approach.

▣ Ans. :

- In supervised approach the complexity of WSD system is handled by machine learning techniques but uses data annotation as well which is the process of adding metadata in the form of tags to the data.
- The words present in a particular corpus are first manually disambiguated.
- Various features are extracted for these words.

- Then the supervised system is designed which has a machine learning classifier or trained on these features.
- The advantage of this system is user can form and update rules in the form of features and also generate training data.
- Many researchers have explored various techniques to address WSD problem starting from machine learning extended to the use of decision lists which uses rich set of features in machine learning framework.
- The researchers further enhanced these features to incorporate different context levels i.e. sentences, paragraphs, micro context etc.

Q.12 What are the different classifiers used in supervised WSD ?

Ans. :

- The popular and efficient classifiers include
- 1. Support Vector Machines (SVMs) :**
 - Support Vector Machine (SVM) model follows supervised machine learning approach.
 - It performs well on limited amount of data.
- 2. Maximum Entropy (Max Ent) Classifier**
 - It is a probabilistic classifier which belongs to class of exponential model.

Q.13 Explain how unsupervised approach can address various issues posed by supervised approach in the problem of word sense disambiguation.

Ans. :

- In case of supervised approach the Word Sense Disambiguation (WSD), process is obstructed by lack of labelled training data for training a classifier to handle every sense of each word in a particular language.
- This issue can be resolved by following solutions :
 1. There can be several instances or occurrences of a word. Based on these instances the clusters are formed covering examples of words of a certain sense. This is called as sense induction through clustering.
 2. Create a matrix for identification of nearness of a given instance of the word, with sets of known senses of a word. Based on this select the closest sense to that instance.
 3. Consider some examples of words of certain sense. Use these seeds, perform iterations on these seeds and enhance the clusters.

Q.14 What are different algorithms which use distance measures to identify senses ?

OR Explain different algorithms of word sense disambiguation by unsupervised approach.

Ans. :

- There are many algorithms designed by various researchers which uses distance measure to identify cases.
- Some of the algorithms are explained below :
 1. Rada proposed an algorithm which uses a matrix to compute shortest distance between two sense pairs in WordNet.

For example, IS-A in WordNet which assumes that different co-occurring words may posses the senses which minimizes the distance in semantic network.

2. Resnick proposed the measure of semantic similarity known as information content in IS-A taxonomy.
3. Agirre and Rigau proposed the approach which extended information content measure to a measure called as conceptual density as shown in Fig. Q.14.1.

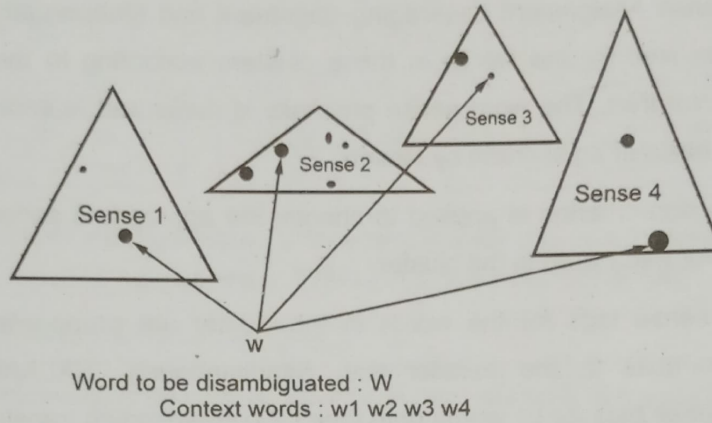


Fig. Q.14.1 : Conceptual density

As shown in Fig. Q.14.1 sense 2 has highest conceptual density so it is chosen as a correct sense.

Conceptual density can be computed as,

$$(D, C_e, m) = \frac{\sum_{i=0}^{M-1} \text{Hyponyms}^{i-0.20}}{\text{Descendants}_c}$$

4. Resnik proposed a measure for computation of sense of a word based on predicate-argument statistics.

A_R is considered to be a selectional association of predicate P to concept C with argument R. A_R is defined as,

$$A_R(P, C) = \frac{1}{S_R(P)} P(C|P) \log \frac{P(C|P)}{P(C)}$$

For noun n, arguments relation R and $\{S_1, S_2, \dots, S_R\}$ as possible senses the

$$C_i = \{C \mid C \text{ is an ancestor of } S_i\}$$

$$a_i = \max_{c \in C_i} A_R(P, C)$$

for $i = 1$ to K where a_i is score of sense S_i

5. Leacock, Miller and Chodorow proposed an algorithm which used corpus statistics and WordNet relations.

Q.15 Explain the algorithm based on crosslinguistic information.

OR Explain SALAAM algorithm.

▣ **Ans. :**

- The algorithm is an unsupervised algorithm based on crosslinguistic information.

- Along with word sense disambiguation it is also helpful in language translation system, as it also figures out-line difference in senses that required translating to other languages.
- The algorithm is explained in Fig. Q.15.1.

1. L1 words that translate into the same L2 word are grouped into clusters.
2. SALAAM (Sense Assignment Leveraging Alignment and Multilinguality) identifies the appropriate senses for the words in those clusters according to the words senses' proximity in WordNet. The word sense proximity is measured in information theoretic terms on the basis of an algorithm by Resnik .
3. A sense selection criterion is applied to choose the appropriate sense label or set of sense labels for each word in the cluster.
4. The chosen sense tags for the words in the cluster are propagated back to their respective contexts in the parallel text. Simultaneously, SALAAM projects the propagated sense tags for L1 words onto their L2 corresponding translations.

Fig. Q.15.1 : SALAAM algorithm for creating training using parallel English-to-Arabic machine translations

Q.16 Explain how WSD can be done using semisupervised approach by Yarowsky algorithm.

OR Explain now Yarowsky algorithm addresses WSD problem.

■ **Ans. :**

- In case of semisupervised approach the process of WSD starts from a small example.
- A classifier is used and an algorithm iteratively identify more training examples.
- Due to this labeled data is generated which is added to training data for more accurate predictive analysis.
- One of such semi-supervised algorithm is Yarowsky algorithm.
- The algorithm takes into consideration two facts :

1. One sense per collocation :

To determine sense of a word, nearly occurring words play an important role based on their relative distance, order and syntactic relationship.

2. One sense per discourse :

In a particular document or for a given discourse a particular word is referred with a same sense.

- The algorithm starts with large untagged corpus. It then identifies the examples of polysemous words (words or phrases having multiple meanings) and stores are relevant sentences as lines.
- To understand the working lets consider the example. Consider a target word plant based on manual choice of possible senses a seed set can be formed for example in our case seed set for word plant can be {life, manufacturing}. The first word relates to flora and second word to industry classified according to rules.

- Consider U as a set of unlabelled examples, as shown in Fig. Q.16.1.
- In the next step of instances of two senses l and m according to the seed set are added to the set as shown in figure.
- This will be done iteratively till no new examples are found.

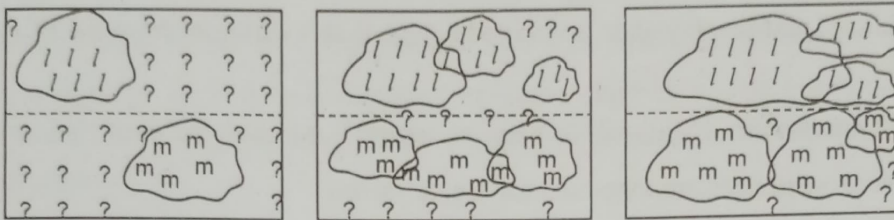


Fig. Q.16.1 : Yarowsky's algorithm : with every iteration new examples are labelled according to seed set {life (l)}, manufacturing (m)}

- The algorithm can be written as shown in Fig. Q.16.2.

Step 1 : In a sufficiently large corpus, identify all the instances of a particular polysemous word that needs to be disambiguated, storing its context alongside.

Step 2 : Identify a small set of instances that are strongly representative of one of the senses of the word. This can either be done in a completely unsupervised fashion by identifying collocations that give a strong indication of the sense usage for the word under consideration or by manually tagging a small portion of the data. In this example, we assume a polysemous word with only two senses, but this algorithm can be extended to n senses.

Step 3 :

Step 3a : Train a supervised classifier on this set of examples.

Step 3b : Using these classifiers, classify the remaining instances of the word in the corpus and select those that are classified above a certain level of confidence.

Step 3c : Filter out the possible misclassifications using one sense per discourse constraint, and identify possible new collocations to be added to the list of seed collocations.

Step 3d : Repeat step 3 iteratively, thereby slowly shrinking the residual.

Step 4 : Stop. At some point, a small, stable residual will remain.

Step 5 : The trained classifier can now be used to classify new data, and that in turn can be used to annotate the original corpus with sense tags and probabilities.

Fig. Q.16.2 : The Yarowsky algorithm

Q.17 Explain the Mihalcea and Moldovan algorithm.

OR Explain the algorithm of for generating examples for words tagged with particular senses.

▣ Ans. :

- The system based on Mihalcea and Moldovan algorithm addresses the problem of lack of hand canot data to train supervised system.

- It works on example of a particular sense from a large corpora in WordNet.
- The algorithm works as follows :

1. Extract : In wikipedia the sentences may contain words as simple links.

For example : `[[bar]]` or a piped link. Piped link is an internal link with wikitext that creates hyperlinked (underlink, clickable) text displayed on a wikipedia page that is different from the title of the page to which the text links.

For example : `[[Train station | Station]]` displays as station but link to the train station wikipedia article.

In this steps all such sentences are extracted for a given word.

- 2. Filter :** If the link is pointing to disambiguation page, the further information is needed to disambiguate the word so filter all such links. If link to disambiguation page is not present then consider the word as label. For piped links string before pipe is a label.
- 3. Collect :** Map all these labels to possible WordNet senses. If they map to various categories of WordNet they provide sense disambiguated data for training.

The algorithm is described in Fig. Q.17.1.

Step 1 : Preprocessing

- For each sense of a word *W*, determine the synsets of WordNet in which it appears. For each such synset, determine monosemous words included in that synset. Parse the gloss definition attached to each synset.

Step 2 : Search

- Form search phrases using the following procedures in order of preference
 1. If they exist, extract monosemous synonyms from the synsets selected in step 1.
 2. Select each of the unambiguous parsed constituents in the gloss as a search phrase.
 3. After parsing the gloss, replace all stop-words with a NEAR operator and create a query from the words in the current synset. For example, if the synset for *produce* #6 is *grow, raise, farm, produce*, and the gloss is *cultivate by growing*, then the query will look like : *cultivate NEAR growing AND (grow OR raise OR farm OR produce)*.
 4. Use only the head phrase combined by words in the synset using the AND operator. For example, if the definition for *company* #5 is *band of people* and its synset is *(party, company)*, then the query becomes : *band of people AND (party OR company)*.
- Search the Internet with the phrases determined in the previous step and gather matching documents.
- From these documents, extract the sentences containing these words.

Step 3 : Postprocessing

- Keep only those sentences in which the word under consideration belongs to the same part of speech as the selected sense, and delete the others.

Fig. Q.17.1 : Mihalcea and Moldovan algorithm for generating examples for words tagged particular senses by querying a very large corpus

Q.18 What are the different software programs available for word sense disambiguation.

▣ **Ans. :**

Following are some of the software programs explored by the researchers for word sense disambiguation.

1. IMS (It Makes Sense) :

- It is a supervised English all words WSD system. The researchers Z_{h1} Zhong and Hwee TOU N_g have developed this system by making use of classifier linear support vector machines and use of various knowledge based features. This system facilitate users to integrate different preprocessing tools, classifiers and additional features.
- It can be accessed through
<http://nlp.comp.nus.edu.sg/software>.

2. WordNet - Similarity - 2.05 :

- WordNet similarity - 2.05 is a perl module which implements a variety of semantic similarity and relatedness measures based on information found in the lexical database WordNet.
- It incorporate the measures of Resnik, Lin, Jiang-conrath, Leacock - Chodorow, Hirst-St.onge, Wu-Pdlmer, Banerjee - Pederson and Patwardhan - Pedersen.
- The input to the modules is pair of words and the output is a numeric value indicating their degree of similarity or relatedness.
- It can be accessed through
<http://search.cpan.org/dist/WordNet-Similarity>

3. WikiRelate :

- WikiRelate is a system developed by Michael struck and simone Paolo Ponzetto.
- It is used for computing relatedness of the words in a structured way and has more coverage than WordNet.
- The dataset used is wikipedia for computation of semantic relatedness and it is compared to WordNet on various benchmarking datasets.
- It can be accessed through :
<http://www.h-its.org/english/research/nlp/download/wiki>

Fill in the Blanks for Mid Term Exams

- Q.1** To extract the meaning out of syntactically arranged sentences is the objective of _____ parsing.
- Q.2** _____ is analysing a text into logical syntactic components and semantics is study of meaning.
- Q.3** The meaningful parts are identified in the text and they are transformed into suitable _____ for higher level task accomplishment.
- Q.4** Semantic parsing is about identifying the _____.
- Q.5** Semantic interpretation facilitates the joining of different components which define _____ representation of the text.
- Q.6** Semantic parsing is the part of _____ interpretation.

UNIT - IV

4

Predicate Argument Structure

Important Points to Remember

- In linguistics, a **corpus** (Plural corpora) or **text corpus** is a language resource consisting of a large and structured set of texts (now a days usually electronically stored and processed). In corpus linguistics, they are used to do statistical analysis and hypothesis testing, checking occurrences or validating linguistic rules within a specific language territory.
- An **annotation** is extra information associated with a particular point in a document or other piece of information. It can be a note that includes a comment or explanation.

☐ Predicate :

- NLP predicates are the words that we use to express our senses like Visual, auditory, kinaesthetic, ADor self-talk (also called labeling system), olfactory(what you smell) and gustatory (what you taste) to the outside world.
- These are also referred to as process words or sense-specific words within NLP. These words tell us in which preferred representation system a person is.

☐ Argument :

- In linguistics, an argument is an expression that helps complete the meaning of a predicate.
- A predicate and its arguments form a predicate-argument structure.
- The discussion of predicates and arguments is associated most with (content) verbs and noun phrases (NPs).
- For e.g. : Devika likes oranges, in this sentence two arguments are there, the first noun (phrase) is the subject argument and the second the object argument.

4.1 Predicate Argument Structure

Q.1 What are the resources to address predicate argument structure ?

OR What are the approaches for converting linguistic insights into features ?

☐ Ans. :

- **FrameNet** and **PropBank** are two major semantically tagged corpora.
- They have evolved from rule based approaches to data oriented approaches.
- They transform linguistic understanding to features other than rules.

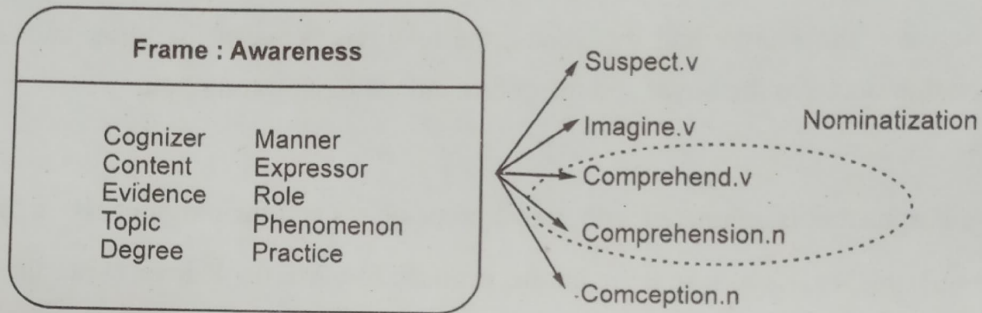
- Machine learning techniques can use these features to train the model which can be used to automatically tag the semantic information which is present in these resources.
- **FrameNet**
 - **FrameNet** is a project housed at the international computer science institute in Berkeley, California which produces an electronic resources based on a theory of meaning called frame semantics.
 - **Frame semantics** is a theory of linguistic meaning developed by Charles J. Fillmore that extends his earlier case grammar. It relates linguistic semantics to encyclopedic knowledge.
 - The basic idea is that one cannot understand the meaning of a single word without access to all the essential knowledge that relates to that word.
 - For example, one would not be able to understand the word “sell” without knowing anything about the situation of commercial transfer, which also involves, among other things, a seller, a buyer, good, money, the relation between the money and the goods, the relations between the seller and the goods and the money, the relation between the buyer and the goods and the money and so on.
- **PropBank**
 - PropBank is a corpus that is annotated with verbal propositions and their arguments - a “proposition bank”.
 - Although “PropBank” refers to a specific corpus produced by Martha Palmer et al., the term PropBank is also coming to be used as a common noun referring to any corpus that has been annotated with propositions and their arguments.

Q.2 Explain FrameNet.

▣ Ans. :

- **FrameNet** is a project developed at the international computer science Institute in Berkeley, California by Charles J. Fillmore.
- It is based on the concept of theory of meaning called frame semantics.
- A semantic frame is a conceptual structure describing an event, relation or object and the participants in it.
- For e.g. Consider the example: “Atul sold a house to Anil” describes the same fact as “Anil bought house from Atul”. This same basic situation is known as a semantic frame.
- The FrameNet lexical database contains over 1, 200 semantic *frames*, 13000 *lexical units* (a pairing of a word with a meaning; polysemous words are represented by several *lexical units*) and 202,000 example sentences.
- The goal of FrameNet project is to facilitate automatic semantic RoleLabeling.
- It contains semantic annotation of predicates in English and also contain tagged sentences from British National Corpus (BNC).
- The process involves :
 - Identification of frames invoked by predicates and creation of frame specific roles called as frame elements.
 - Identification of predicates which instantiated the frame and label the sentences for those predicates.

- The labeling process involves :
 - Identification of frame instantiated and invoked by predicate lemma.
 - Identification of semantic arguments for that instance.
 - Tagging these arguments with predetermined set of frame elements for that particular frame.
- Combination of predicate lemma and the frame invoked by its instance is known as Lexical Unit (LU).
- For e.g. : Consider The DRIVING frame, for example, specifies a DRIVER (a principal MOVER), a VEHICLE (a particularization of the MEANS element) and potentially CARGO or RIDER as secondary movers. In this frame, and DRIVER initiates and controls the movement of the VEHICLE. For most verbs in this frame, DRIVER or VEHICLE can be realized as subject; VEHICLE, RIDER or CARGO can appear as direct objects; and PATH and VEHICLE can appear as oblique complements.
- For example in Fig. Q.2.1 Frame is : AWARENESS, verb predicate : Believe, noun predicate : comprehension.



1. [Cognizer We] [Predicate : verb believe] [Content it is a fair and generous price]
2. No doubts existed as [Cognizer OUR] [Predicate:noun comprehension] [Content of it]

Fig. Q.2.1 : FrameNet example

Q.3 Explain PropBank.

▣ Ans. :

- PropBank is a corpus of text annotated with information about basic semantic propositions created by Martha Palmer.
- It is created by adding predicate-argument relations to the syntactic trees of Penn Treebank.
- In PropBank the arguments of each predicate are annotated with their semantic roles in relation to the predicate.
- PropBank annotation also requires the choice of a sense id (also known as a 'frameset' or 'roleset' id) for each predicate.
- For each verb in every tree which presents the phrase structure of the sentence, a PropBank instance is created which consists of the sense id of the predicate and its arguments labeled with semantic roles.
- It provides consistent argument labels across different syntactic realizations of the same verb.

For e.g. : [ARG0 Devika]broke [ARG1 the window]

[ARG1 The window] broke

- In the above example the arguments of the verbs are labeled as numbered arguments : Arg0, Arg1, Arg2 and so on.
- As shown in Table Q.3.1 the arguments are tagged as core arguments (label type ARGN, where N = 0 to 5) or adjunctive arguments (label type ARGM-X, where X = TMP for temporal, LOC for locative, etc.)

Table Q.3.1 : List of adjunctive arguments in PropBank - ARGMS

Tag	Description	Example
ARGM-LOC	Locative	<i>the museum, in Westborough, Mass</i>
ARGM-TMP	Temporal	<i>now, by next summer</i>
ARGM-MNR	Manner	<i>heavily, clearly, at a rapid rate</i>
ARGM-DIR	Direction	<i>to market, to Bangkok</i>
ARGM-CAU	Cause	<i>In response to the ruling</i>
ARGM-DIS	Discourse	<i>for example, in part, Similarly</i>
ARGM-EXT	Extent	<i>at \$38.375, 50 points</i>
ARGM-PRP	Purpose	<i>to pay for the plant</i>
ARGM-NEG	Negation	<i>not, n't</i>
ARGM-MOD	Modal	<i>can, might, should, will</i>
ARGM-REC	Reciprocals	<i>each other</i>
ARGM-PRD	Secondary predication	<i>to become a teacher</i>
ARGM	Bare ARGM	<i>with a police escort</i>
ARGM-ADV	Adverbials	<i>(none of the above)</i>

- For e.g. Core arguments for predicates operate and author are shown in Table Q.3.2.

Table Q.3.2 : Argument labels associated with the predicate operate.01 (sense : work) and for author.01 (sense : to write or construct) in the PropBank corpus

Predicate	Argument	Description
operate.01		
	ARG0	Agent, operator
	ARG1	Thing operated
	ARG2	Explicit patient (thing operated on)
	ARG3	Explicit argument
	ARG4	Explicit instrument
author.01		
	ARG0	Author, agent
	ARG1	Text authored

Q.4 Explain semantic role labeling using Semantic Role Labeling (SRL) algorithm.

Ans. :

- Gildea and Jurafsky proposed that semantic role labeling is a unsupervised classification problem, in which predicate and arguments of predicate are mapped to a node in a syntax tree for a particular sentence.
- Semantic role labeling consist of three tasks :
 - **Argument identification** : Involves identification of all the parse constituents representing valid semantic arguments of a predicate.
 - **Argument classification** : Involves assigning appropriate argument labels to identified constituents.
 - **Argument identification and classification** : It is a combination of above two tasks i.e. identification of constituents and assigning the label to them.
- The process can be explained through Semantic Role Labeling (SRL) algorithm shown in Fig. Q.4.1.

Procedure : SRL (sentence) returns best *semantic role labeling***Input** : *Syntactic parse of the sentence*

1. *Generate a full syntactic parse of the sentence*
2. *Identify all the predicates*
3. *For all predicate \in sentence do*
4. *Extract a set of features for each node in the tree relative to the predication*
5. *Classify each feature vector using the model created in training*
6. *Select the class of highest scoring classifier*
7. *Return best semantic role labeling*
8. *end for*

Fig. Q.4.1 : The Semantic Role Labeling (SRL) Algorithm**Q.5 Explain the features stated by Gildea and Jurafsky in Phrase Structure Grammar (PSG) for semantic role labeling problem.**

▣ Ans. :

- In semantic role labeling problem apart from frameNet and propBank, certain high quality statistical parsers can be used to generate a phrase structure tree.
- As phrase structure is open for modifications Gildea and Jurafsky used some features to address semantic role labeling problem.
- These features are explained below

1. Path

- Consider the example in Fig. Q.5.1, the path which starts from ARGO. It to predicate operates follows the string.

NP → S → VP → VBZ

OR NP ↑ S ↓ VP ↓ VBZ

Where ↑ = Upward movement , ↓ = Downward movement in tree

o This is known as a Syntactic Path from parse constituent to predicate be classified.

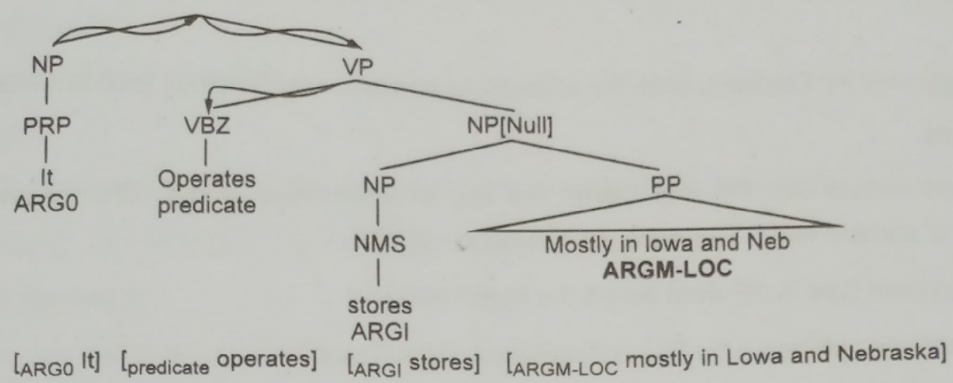


Fig. Q.5.1 : Syntax tree for a sentence illustrating the PropBank tags

2. Predicate

This includes identification of predicate lemma.

3. Phrase type

It identifies the syntactic category.

Ex. NP, PP, S, etc. of the constituent to be labeled.

4. Position

This feature identifies if the phrase is present before or after the predicate.

5. Voice

This feature indicates if the predicate is in active or passive role. For passive voiced predicates $tgrep2^7$ expressions are used on syntax tree.

6. Head word

This feature represents syntactic head of phrase and is calculated using head word table.

7. Subcategorization

- o This expands the predicates parent node in parse tree.
- o For e.g. in Fig. Q.5.1, for predicate operates the subcategorization is given as,
 $VP \rightarrow VBZ - NP$

8. Verb clustering

- o In certain cases it may happen that hand tagged training data is very limited and for any real time test set the predicates may not be present in the training.
- o In these cases the predicate cluster can be formed and it can be used as a feature.
- o For example, the verbs like eat, devour, savour are the object which describe food.
- o The distance function is used for clustering based on the fact that verbs with similar semantics will have similar objects.

Q.6 Explain the features suggested by Surdeanu et al. for semantic role labelling problem.

▣ **Ans. :**

The features suggested by Surdeanu et al. for addressing semantic role labelling problem includes :

1. Content word

As head word feature like PP and SBAR are not so informative, some different rule set is used for identification of content word. The rules are shown in Fig. Q.6.1.

H1 : if phrase type is PP **then** select the rightmost child

Example : phrase = "in Texas," content word = "Texas"

H2 : if phrase type is SBAR **then** select the leftmost sentence (S*) clause

Example : phrase = "that occurred yesterday," content word = "occurred"

H3 : if phrase type is VP **then**

if there is a VP child **then**

select the leftmost VP child

else

select the head word

Example : phrase = "had placed," content word = "placed"

H4 : if phrase type is ADVP **then** select the rightmost child, not IN or TO

Example : phrase = "more than," content word = "more"

H5 : if phrase type is ADJP **then** select the rightmost adjective, verb, noun or ADJP

Example : phrase = "61 years old," content word = "61"

H6 : for all other phrase types select the head word

Example : phrase = "red house," content word = "red"

Fig. Q.6.1 : List of content word heuristics

2. Part of speech of the head word and content word

Part Of Speech (POS) tags can increase the performance the decision tree system.

3. Named entity of the content word

- The roles like ARGM-TMP and ARGM-LOG represent TIME or PLACE named entities.
- This information can be represented as set of binary values.

4. Boolean named entity flags

Named entities like PERSON, PLACE, TIME, DATE etc. are considered as features.

5. Phrasal verbs collections

The frequency of verbs and immediately following preposition is computed and it is used as a feature.

Q.7 Explain the features added by Fleischman, kwon and Hovy to their designed system for semantic role labeling problem.

Ans. : Following features are added :

1. Logical function

- This function take as input external argument object argument and other argument.
- It is computed by some heuristics on syntax tree and is used as feature.

2. Order of frame elements

Relative position of a frame element to other frame elements in a sentence is considered as a feature.

3. Syntactic pattern

Heuristics on phrase type and logic function of the constituent is used as a feature.

4. Previous role

n^{th} previous role assigned by system for current predicate is considered as a feature.

Q.8 Explain Combinatory Categorical Grammar (CCG).

Ans. :

- In argument identification task, path feature is very important.
- But it is difficult to train feature.
- According to the research of Gildea and Hockenmater, due to features from CCG semantic role labeling is improved.
- Consider the example shown in Fig. Q.8.1 which shows CCG of the sentence "London denied plans on Monday".

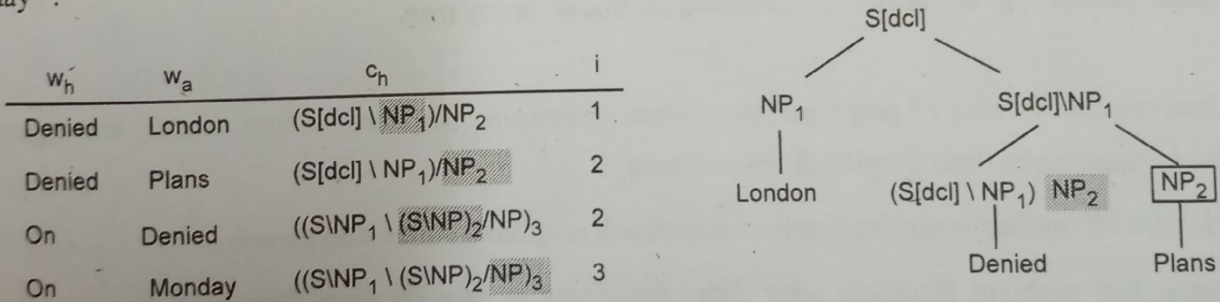


Fig. Q.8.1 : Combinatory category grammar parse

• Three features are proposed :

1. Phrase type : It indicates maximum projection between predicate and dependent word.

2. Categorical path : This feature contain three values

1. Dependent word category.
2. Direction of dependence.
3. Slot in the category filled by dependent word. For example, the path between denied and plans.

3. Tree path : It is a path from dependent word through predicate.

Q.9 Explain how Tree-Adjoining Grammar (TAG) is used for semantic role labeling.

▣ **Ans. :**

- Tree Adjoin Grammar (TAG) is conceptual used by Chen and Rambow and it takes into account the long distance dependencies.
- They used two sets of features to address semantic role labeling problem in predicate argument structure
 1. Surface syntactic features.
 2. Additional features by extracting TAG from Penn Treebank.
- These additional features include :
 1. **Supertag path** : It is a path features derived from TAG instead of PSG.
 2. **Supertag** : Tree frame feature related to predicate or argument.
 3. **Surface syntactic role** : This feature includes syntactic role of the argument.
 4. **Surface sub categorization** : It is a sub categorization frame feature.
 5. **Deep syntactic role** : This feature includes deep syntactic role of an argument. The values are subject and direct object.
 6. **Deep sub categorization** : It is deep sub categorization frame. For example, NPO → NP1 for transitive verb.
 7. **Semantic sub categorization** : In addition to semantic sub categorization frame, semantic role information is used.

Q.10 Explain problem of semantic role labeling on dependency tree.

Ans. :

- It is observed that in case of propBank the system performance depends on how exactly the arguments are annotated according to Penn Tree-bank constituents.
- Correct score of labelled is obtained only if they match propBank annotation exactly.
- To address this problem Hacioglu used dependency tree to convert Penn Tree-bank trees to dependency representation.
- The performance is observed to be increased by SF score point then on phrase structure trees.
- Dependency is established between a word called as head and another called as modifier and a dependency tree is generated.
- Each word can modify at most one other word.
- The headword are classified based on the features shown in Table Q.10.1.

Table Q.10.1 : Features used in the baseline system using Minipar parses

Head word	The word representing the node in the dependency tree.
Head word POS	Part of speech of the head word.
POS path	The path from the predicate to the head word through the dependency tree connecting the part of speech of each node in the tree.
Dependency path	Each word that is connected to the head word has a dependency relationship to the word. These are represented as labels on the arc between the words. This feature comprises the dependencies along the path that connects two words.
Voice	Voice of the predicate.
Position	Whether the node is before or after the predicate.

- In Fig. Q.10.1, we can see the representation of word kick according to phrase structure grammar and minipar parse tree.

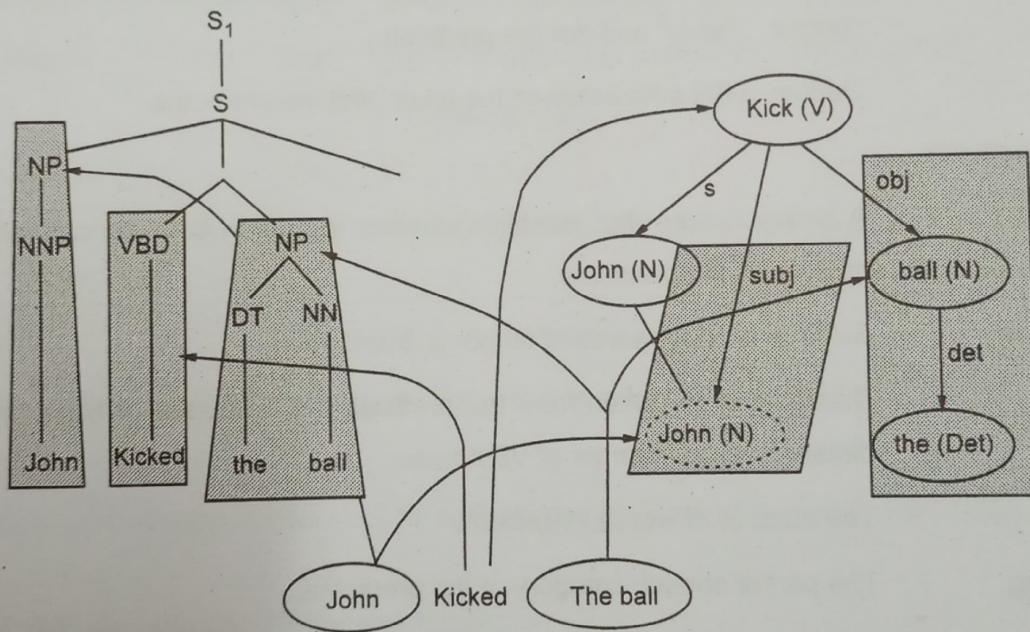


Fig. Q.10.1 : New architecture

Q.11 Explain chunk based approach for semantic roll labeling.

OR Explain base phrase chunks.

▣ **Ans. :**

- It is observed that in case speech data the chunk representation is faster and robust as compared to phrase reordering.
- In chunk based system each base phrase is classified as B (eginning) of a semantic role, I (inside) a semantic role or o(utside) semantic role.
- This is known as IOB representation.

- In this system the input text is first chunked or divided into base phrases by SVM classifier.
- In the second step, second SVM is trained for assignment of semantic labels to chunks.
- The Table Q.11.1 shows the features used by semantic chunker and Fig. Q.11.1 shows the semantic chunker system.

Table Q.11.1 : Features used by chunk-based classifier

Words	Words in the chunk
Predicate lemma	The predicate lemma
POS tags	Part of speech of the words in the chunk.
BP positions	The position of a token in a phrase (BP) using the IOB2 representation [e.g. B-NP, I-NP, O]
Clause tags	The tags that mark token positions is a sentence with respect to clauses.
Named entities	The IOB tags of named entities.
Token position	The position of the phrase with respect to the predicate has three values : "before,," "after" and (for the predicate).
Path	Defines a flat path between the token and the predicate.
Clause bracket patterns	
Clause position	A binary feature that identifies whether the token is inside or outside the clause containing the predicate.
Headword suffixes	Suffixes of head words of length 2, 3 and 4.
Distance	Distance of the token from the predicate as a number of base phrases and the distance as the number of VP chunks.
Length	The number of words in a token.
Predicate POS tag	The part of speech category of the predicate.
Predicate frequency	Frequent or rare using a threshold of 3
Predicate BP context	The chain of BPs centered at the predicate within a window of size $- 2 / + 2$.
Predicate POS context	POS tags of words immediately preceding and following the predicate.
Number of predicates	This is the number of predicates in the sentence.

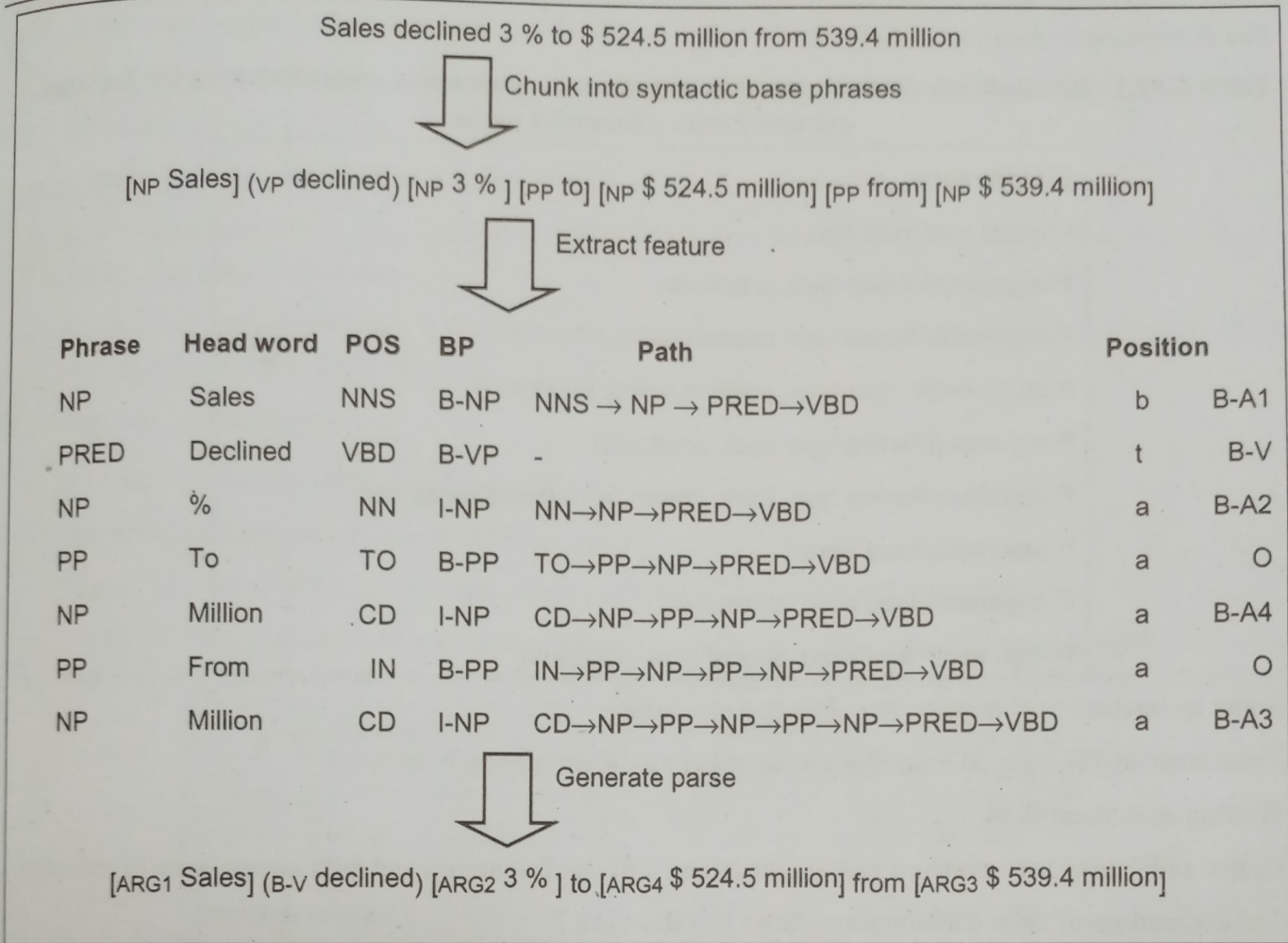


Fig. Q.11.1 : Semantic chunker

Q.12 Explain how machine learning algorithms are used to address semantic role labelling problem.

- Ans. :
- Some of the high performance approaches to address the problem of semantic role labeling are explained below :
 - Variation of SRL algorithm by Gildea and Jurafsky.
 - It consists of two steps :
 - Step 1 :** Calculation of maximum likelihood probabilities to know that a constituent is an argument based on features.
 - P (argument | Path, predicate) and
 - P (argument | Head, predicate)
 - Step 2 :** Assign each constituent having nonzero probability a normalized probability.

- The distribution is shown in Table Q.12.1.

Table Q.12.1 : Distributions used for semantic argument classification, calculated from the features extracted from a Charniak parse

Distributions
P(argument Predicate)
P(argument Phrase type, predicate)
P(argument Phrase type, position, voice)
P(argument Phrase type, position, voice, predicate)
P(argument Phrase type, path, predicate)
P(argument Phrase type, Path, predicate, subcategorization)
P(argument Head word)
P(argument Head word, predicate)
P(argument Head word, phrase type, predicate)

- Some researchers used decision tree classified algorithm.
- Fleischman and Hovy used FrameNet corpus using maximum entropy framework.
- Pradhan et.al. used SVM.
- Gildea and Palmer used posterior probabilities using different feature sets and does interpolation of estimates.
- The comparison of those classifiers are shown in Table Q.12.2.

Table Q.12.2 : Argument classification using same features but different classifiers

Classifier	Accuracy (%)
SVM (Pradhan et.al.) [120]	88
Decision tree (Surdeanu et.al.) [118]	79
Gildea and Palmer [131]	77

Q.13 What are the different strategies to handle the limitation of treating semantic role labelling as a series of independent argument classification.

Ans. : Following are some of the strategies :

1. Disallowing overlaps :

- Every constituent is classified independently.
- It may happen that overlapping constituents are assigned same argument type as shown in the below example.

But [_{ARG0} nobody] [_{Predicate} knows] [_{ARG1} at what level] [_{ARG1} the features and stocks, will open today]

- As overlapping arguments are not allowed in propBank this problem can be solved by choosing the argument for which SVM has highest confidence based on classification probabilities and labeling others as null.

□ 2. Argument sequence information :

- This assumes that to increase performance to argument tagger predicate can instantiate a certain set of arguments.
- Argument ordering information is retained and predicate is considered as a part of argument.
- This is accomplished by first converting raw, SVM scores to probabilities and then crediting argument lattice using a best hypothesis.
- In the next step viterbi search is performed according to the probabilities assigned by sigmoid and maximum likelihood path is obtained.

□ 3. Feature performance :

- It is observed that each feature and useful for every task.
- Then efficiency depends on classification paradigm.
- Table Q.13.1 shows effect of each feature on argument classification and identification.

Table Q.13.1 : Effect of each feature on the argument classification task and argument identification task when added to the baseline system. An asterix indicates that the improvement is statistically significant

Features	Argument classification	Argument identification		
		P	R	F ₁
Baseline [120]	87.9	93.7	88.9	91.3
+ Named entities	88.1	93.3	88.9	91.0
+ Head POS	*88.6	94.4	90.1	*92.2
+Verb cluster	88.1	94.1	89.0	91.5
+ Partial path	88.2	93.3	88.9	91.1
+ Verb sense	88.1	93.7	89.5	91.5
+ Noun head PP (only POS)	*88.6	94.4	90.0	*92.2
+ Noun head PP(only head)	*89.8	94.0	89.4	91.7
+ Noun head PP (both)	*89.9	94.7	90.5	*92.6
+ First word in constituent	*89.0	94.4	91.1	*92.7
+ Last word in constituent	*89.4	93.8	89.4	91.6
+ First POS in constituent	88.4	94.4	90.6	*92.5
+ Last POS in constituent	88.3	93.6	89.1	91.3
+ Ordinal constituent pos. concat.	87.7	93.7	89.2	91.4
+ Const. tree distance	88.0	93.7	89.5	91.5
+ Parent constituent	87.9	94.2	90.2	*92.2
+ Parent head	85.8	94.2	90.5	*92.3
+ Parent head POS	*88.5	94.3	90.3	*92.3
+ Right sibling constituent	87.9	94.0	89.9	91.9
+ Right sibling head	87.9	94.4	89.9	*92.1
+ Right sibling head POS	88.1	94.1	89.9	92.0

Features	Argument classification	Argument identification		
		P	R	F ₁
	A			
+ Left sibling constituent	*88.6	93.6	89.6	91.6
+ Left sibling head	86.9	93.9	86.1	89.9
+ Left sibling head POS	*88.8	93.5	89.3	91.4
+ Temporal cue words	*88.6	-	-	-
+ Dynamic class context	88.4	-	-	-

4. Feature salience :

- Relative contribution of different feature sets as shown in Table Q.13.2 is a contributing factor in system performance.

Table Q.13.2 : Performance of various feature combinations on the task of argument classification

Features	Accuracy
All features [120]	91.0
All except path	90.8
All except phrase type	90.8
All except HW and HW-POS	90.7
All except all phrases	*83.6
All except predicate	*82.4
All except HW and FW and LW info.	*75.1
Only path and predicate	74.4
Only path and phrase type	47.2
Only head word	37.7
Only path	28.0

- The system performance decreases one feature is left at a time.
- In Table Q.13.3, the importance of features on argument identification is shown.

Table Q.13.3 : Performance of various feature combinations on the task of argument identification

Features	P	R	F ₁
All features [120]	95.2	92.5	93.8
All except HW	95.1	92.3	93.7
All except predicate	94.5	91.9	93.2
All except HW and FW and LW info.	91.8	88.5	*90.1
All except path and partial path	88.4	88.9	*88.6
Only path and HW	88.5	84.3	86.3
Only path and predicate	89.3	81.2	85.1

5. Feature selection :

- Feature selection strategy is important as it plays different roles in argument classification and argument identification.
- If named entity features are added argument identification task is deteriorated whereas argument classification task is improved.
- Different feature selection strategies can be adapted. Some of them are :
 - Leave one feature from full set considering that features are independent and check the performance. Based on the effect either keep the feature or leave it.
 - Use of SVM can be done. But the drawback is SVM output distances instead of probabilities so it is difficult to compare them cross classifiers as different features are used to train classifiers.
 - Solution to this is conversion of SVM scores to probabilities.
 - Also Pool Adjustment Violators (PAV) algorithm is purposed by foster and stine to convert classifier scores to probabilities.

Q.14 Explain how size of training data effects performance of classifier in supervised learning method.

Ans. :

- As per the research done by pradhan et. al the effect of the amount of training data identification and classification of arguments is shown in Fig. Q.14.1.
- The topmost curve in the figure indicates F_1 score variation for argument identification task.
- Third curve indicates F_1 score on argument identification and classification.
- After 10,000 examples a plateau is reached indicating that only tagging more amount of data will not help to achieve the performance.

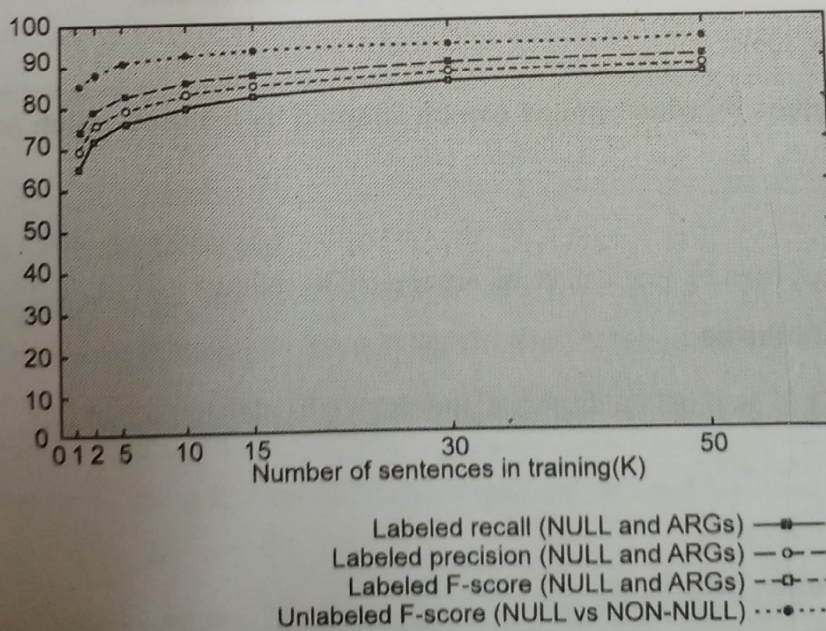


Fig. Q.14.1 : Learning curve for the tasks of identifying and classifying arguments using treebank parses

Q.15 How to overcome parsing errors in predicate argument structure.

▣ Ans. :

- It has been observed by the researches pradhan et.al. that argument identification is a problem which affects overall systems performance.
- These errors are caused by the fact that
 1. Failure of syntactic analyzer in providing constituents mapping to correct arguments.
 2. Failure of a system to identify the constituents corresponding to semantic roles.
- These errors are overcome by following two techniques
 1. Combining parses from different syntactic representations.
 2. Use of best parses.

Q.16 Explain how semantic role labeling is done for nominal predicates or normalizations.

Ans. :

- In a sentence consists of various words consisting of different part of speech tags.
- So it is along with verbs it is also necessary to identify arguments of predicates like nominal, adjectival and pre-positional.
- One way to handle nouns is nominalization i.e. converting a verb into abstract noun.
- Consider the example of nominalization as shown in Fig. Q.16.1.

She complained **about the attack**

She ~~made~~ an official complaint **about the attack**

Nominalized sentence

John walked **around the university**

John ~~took~~ a walk **around the university**

Nominalized sentence

Fig. Q.16.1 : Example of nominalization

- The verbs are make and took.

Q.17 Explain the techniques by which nouns can be adapted by features which are originally devised for verbs.

▣ Ans. :

Some of the features proposed by Pradhan et. al. are limited as follows :

▣ 1. **Intervening verb features :**

- To realize arguments of nominal predicates supporting verbs play important role.
- Three classes being used are :
 1. Verbs by being
 2. Light verbs like make, take, have.
 3. Verbs with POS starting with VB.
- Three features added are :
 1. A binary feature which shows presence of verb between predicates and arguments.

2. Actual word as a feature.
3. Path from constituent to verb in a tree.

For example : [Speaker Amit] makes general [_{predicate} assertions] [Topic about marriage]

□ 2. **Predicate NP expansion rule :**

- Lower most NP is located and expansion rule is applied to it.
- This features clusters noun phrases with similar internal structure and helps in finding modifiers.

□ 3. **Is predicate plural :**

This features tells if predicate is singular or plural.

□ 4. **Genitives in constituent :**

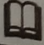
- Genitive words are the ones with POS PRP, PRP\$ or WP\$.)
- This features takes value as true if there is genitive word present.

For example [_{Speaker} Burma's] [_{Phenomenon} oil] [_{Predicate} Search] hits virgin forests.

Q.18 Explain various issues in semantic role labelling faced by the languages other than english.

□ **Ans. :** Following are some of the issues, advantages and disadvantages of the systems. Performing semantic role labelling of the languages other than english :

1. Some language specific features like predicate frame feature for chinese language are beneficial for english processing as well.
2. There are certain features which are present only in some languages.
For example, there are no delimiters to separate the words in chinese language. So to do word segmentation more complex system is needed, where as english language this task become easy.
3. The language like chinese is morphological poor language so predicates and arguments posses close connection as verbs, nouns and adjectives have similarity.
4. Chinese has more verb types than english. So in corpus of same size as english the instances of chinese verbs are less posing the issue of data sparsely.
5. In chinese language syntactic parsing yields less performance than english so shallow parsers give more promising results than full syntactic parsing.
6. In the arabic language POS categories are more than english or chinese, so it is challenging to handle all these by semantic role labeling system and still the rich structure is not explored.

 **4.2 Meaning Representation Systems**

Q.19 What is deep semantic parsing ? Explain various resources used for deep semantic passing.

□ **Ans. :**

- The biggest challenge faced by natural language processing system is the ambiguity which exists at every phase.
- It is challenging task to take input in natural language, remove ambiguity, use word knowledge and understand the context and make machines understand the language effortlessly like humans.

- Still the research is going on resolving various ambiguities of the language and till now the systems are developed for working in specific domains, instead of the generalized models.
- This is known as deep semantic parsing.
- Various resource which are developed in this area include :

□ 1. ATIS :

- **Air Travel Information System (ATIS)** is a first system which transformed natural language input information for decision making by and application.
- It takes as an input user speech query about flight information in restricted vocabulary.
- It converts it into SQL query and retrieve information from flight database.
- Semantic information is encoded from hierarchical frame representation.
- The training corpus consists of 774 sceneries, 137 subjects, over 7300 utterances.
- 2900 utterances are categorized are represented with reference answers.
- 600 are tree-banked.
- The example user query in shown in Fig. Q.19.1.

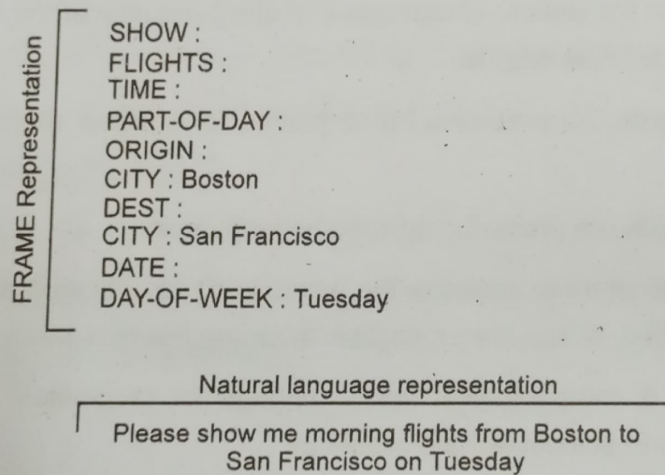


Fig. Q.19.1 : A sample user query and its frame representation in the ATIS program

□ 2. Communicator :

- It is an extension of ATIS system.
- ATIS system was designed for user initiated dialog i.e. answers were provided by machine for user questions.
- Communication system uses mixed initiative dialog. i.e. it is an interactive system in which the communications happens between system and user. Based on real time information provided by system the user can plan the itinerary.
- Many dialogs were collected by the system and stored in linguistic data consortium.
- The data was also collected Carnegie Mellon University.
- About million words and 1600 dialogs are annotated.

□ 3. Geo-Query :

- The system is designed for a domain of US geography.
- It uses Natural Language Interface (NLI) to geographic database geobase.
- The geographic information like population, neighboring states, major rivers and cities are stored in the relational database in the form of 800 prolog facts.
- Sample queries take the form as
answer (C, C capital (S, C) largest (P, (State (S), Population (S, P))))
- Such queries from Geo-Query corpus and is translated in Japanese, Spanish and Turkish.

□ 4. Robocup : CLang :

- This resource use robotic soccer as domain.
- CLang formal language is used, which takes input from team coach and generate output in the form of behaviors as if - then rules.
- For example, consider the advice by coach :
If the ball is in our penalty area, all our players except player 4 should stay in our half
- It is converted to rule as
((POS) (Penalty-area our)) (do (Player-except our 4) (POS (half our)))

Q.20 How natural language is mapped to the meaning representation systems (SQL query, prolog or domain-specific query).

□ Ans. : The different ways by which is a natural language input is mapped to domain specific meaning representation system are follows :

□ 1. Rule based :

- Rule based semantic parsing system give better results for ATIS and communicator systems.
- In these the speech recognition errors are handled by an interpreter.
- Meaning units in the sentence are parsed to semantic structure, based on the fact that syntactic explanation is complex than semantic information.
- Due to the dynamic variations in the spontaneous speech, pauses, stutters word order is considered to be less important which results in scattering of meaning units and does not cater to the order required by syntactic parser.
- Phoenix is such a system designed by word's based on recursive transition networks and handcrafted grammar.
- The system takes into account hierarchical frame structure and the values of these frames are adjusted with new piece of information.
- The error rate observed is,
Spontaneous speech : Input 13.2 % with word error rate 4.4 %.
Transcript input : 4.4 %

2. Superoised :

- Rule based systems have following shortcomings :
 1. Effects for creation of rules.
 2. They are time consuming which affects development of systems.
 3. Difficult to maintain.
 4. Less scalable in case of complex problem.
- To overcome these difficulties statistical models with hand annotated data can be used.
- But if hand annotated data is not present then the system fails.
- To address this problem Schwartz et. al. devised end to end superoised statistical learning system for ATIS.
- The system has four components :
 1. Semantic parse
 2. Semantic frame
 3. Discourse
 4. Backend
- Human-in-the-loop corrective approach is used for training to obtain more data for supervision.

Fill in the Blanks for Mid Term Exam

- Q.1 NLP _____ are the words that we use to express our senses like Visual, auditory, etc.
- Q.2 Predicates are also referred to as _____ words or _____ words within NLP.
- Q.3 In linguistics, an _____ is an expression that helps complete the meaning of a predicate.
- Q.4 A predicate and its arguments form a _____ structure.
- Q.5 **FrameNet** and **PropBank** are two major _____ tagged corpora.
- Q.6 _____ is a project housed at the international computer science institute in Berkeley, California which produces an electronic resources based on a theory of meaning called frame semantics.
- Q.7 _____ is a corpus that is annotated with verbal propositions and their arguments - a "proposition bank".
- Q.8 A _____ is a conceptual structure describing an event, relation or object and the participants in it.
- Q.9 The goal of FrameNet project is to facilitate automatic semantic _____.
- Q.10 Identification of frames invoked by predicates and creation of frame specific roles called as _____.
- Q.11 Combination of predicate lemma and the frame invoked by its instance is known as _____.
- Q.12 _____ is a corpus of text annotated with information about basic semantic propositions created by Martha Palmer.
- Q.13 _____ strategy is important as it plays different roles in argument classification and argument identification.

Multiple Choice Questions for Mid Term Exam

- Q.1 In NLP predicates are the words which express _____.
- a human senses
- b text
- c language
- d syntax

Discourse Processing and Language Modeling

Part : I - Discourse Processing

5.1 : Cohesion

Q.1 What is cohesion ?

▣ Ans. :

- Typically a high level text documents like for ex. Academic articles are divided into different sub sections like Abstract, Introduction, Methodology, Result and Conclusion.
- Automatic detection of all these types is a difficult problem.
- To address this problem the algorithms for discourse segmentation are used.
- The unsupervised discourse segmentation algorithms are based on concept of cohesion.
- Cohesion is linking of different textual units by using linguistic devices.
- Lexical cohesion is the cohesion which exists in two units which is indicated by relations between words in those units.

For ex. : Consider the sentence

Peel, core and slice – the pears and apples.

Add the fruit to skilled.

- In this sentence lexical cohesion between these two sentences is shown by hypernym relation between fruit and words pears and apples.

5.2 : Reference Resolution

Q.2 What is reference resolution ?

▣ Ans. :

- To understand the concept of reference resolution, lets consider following example.
“Lakshmi is studying in 10th standard. She is a very good singer and participated in many music programs. The performance of this 16 years old is also excellent in academics.”
- In the above passage there is mention of one person named Lakshmi.
- The linguistic expressions like her, she are used to denote an individual is known as reference.
- Reference resolution is a task to determine what entities are referred to by which linguistic expressions.
- Referring expression (for ex. she) is a haliral language expression used to perform reference.

- The referred entity is called as reference (for ex. Lakshmi).
- Reference to an entity which is previously introduced in discourse is called anaphora and the referring expression is called as anaphoric. For ex. Pronoun, she and 16 years old are anaphoric.
- Two referring expressions used to refer same entity are said to corefer.
- So typically the task of reference resolution involve two tasks :
 1. Coreference resolution
 2. Pronominal anaphora resolution
- Coreference resolution is the task to find out referring expressions referring to same entity.
- Pronominal anaphora resolution is the task to find out antecedent for single pronoun.
For ex. : antecedent of she is Lakshmi. We need to find out the given a pronoun she, its antecedent is Lakshmi.
- Pronominal anaphora resolution can be considered as subtask of coreference resolution.
- There are various algorithms which can be used for this purpose. Some of them are :
 1. Hobbs algorithm
 2. Centering algorithm
 3. Log linear algorithm

Part : II - Language Modelling

5.3 : N - Gram Models

Q.3 What are n gram models?

▣ Ans. :

- To understand concept of n gram models consider the example sentence.
- Please turn your homework.....
- We want to predict the next word which can be 'in', 'over' but definitely it will not the word 'the'. This is known as word prediction and can be done using probabilistic models called as n-gram models.
- In n-gram model from sequence of n-token words the next word is predicted from previous n-1 words.
- It is difficult to compute probability of any word sequence w.
- It can be computed by decomposing it based on chain rule of probability as :

$$P(w) = P(w_1, \dots, w_n) = P(w_1) \prod_{i=2}^n P(w_i | w_{i-1}, w_{i-2}, \dots, w_2, w_1)$$

But as individual product terms also cannot be computed directly n-gram approximation will be useful.

By considering the assumption of history equivalence class that n-1 are useful for predicting a given word, n-gram model can be defined as :

$$P(w) \approx \prod_{i=1}^n P(w_i | w_{i-1}, \dots, w_{i-n+1})$$

- This is based on markov assumption that current word only depends on $n - 1$ preceding words and independent of all the other given words.
- So n -gram model is also called as $(n - 1)$ - th order Markov model.
- Based on length of n the models can be formalized as : for $n = 2 \rightarrow$ biagram \rightarrow considering two word sequence of words ex. "please turn" or "turn your" for $n=3 \rightarrow$ trigram \rightarrow considering three word sequence of words ex: "please turn your" or "turn your homework" and so on for $n = 4,5$, etc.

5.4 : Language Model Evaluation

Q.4 Write a note on language model evaluation.

Ans. :

- To evaluate a performance of language model the best way is to include it in an application and check the performance of that application this is called as extrinsic evaluation.
- But as it is expensive at times, another way of evaluation i.e. intrinsic evaluation is used.
- In intrinsic evaluation a metric is used to quickly evaluate the improvements in language model.
- Two criterias used for intrinsic evaluation of language model are.

1. Coverage rate

- It measures percentage of n -grams in test set
- Sometimes there are cases where some unknown words appear they are called as Out Of Vocabulary (OOV) words which cannot be handled by this type.

2. Perplexity

- It considers the fact that among two probabilistic models the model which fits the test data is the better one.

5.5 : Parameter Estimation

Q.5 Explain maximum likelihood estimation and smoothing.

Ans. :

- Generally n -gram probabilities are estimated by combining maximum likelihood criterion with parameter smoothing.
- The maximum likelihood estimate can be obtained as

$$P(w_i | w_{i-1}, w_{i-2}) = \frac{C(w_i, w_{i-1}, w_{i-2})}{C(w_{i-1}, w_{i-2})}$$

Where $C(w_i, w_{i-1}, w_{i-2})$ is count of trigram w_{i-2}, w_{i-1}, w_i in training data.

- Smoothing is the process of flattening the peaks in n -gram probability distribution by redistributing probability mass. Also zero estimates are replaced by some small nonzero values.
- One of the common smoothing technique is called as back off.
- It splits n -grams whose count in training data fall below predetermined threshold T and also whose count exceed the threshold.

- The backed off probability P_{BO} for w_i given w_{i-1}, w_{i-2} is computed as

$$P_{BO}(w_i | w_{i-1}, w_{i-2}) = \begin{cases} d_c P(w_i | w_{i-1}, w_{i-2}) & \text{if } C \geq T \\ d(w_{i-1}, w_{i-2}) P_{BO}(w_i | w_{i-1}) & \text{otherwise} \end{cases}$$

Where $C \rightarrow$ count of (w_i, w_{i-1}, w_{i-2})

$d_c \rightarrow$ discounting factor applied to higher order distribution.

Q.6 Explain Bayesian parameter estimation method.

▣ Ans. :

- It is a parameter estimation method in which set of parameters of a model are considered as random variable which is governed by previous statistical distribution posterior distribution,
- $P(\theta|S)$ is given using Bay's rule as

$$P(\theta|S) = \frac{P(S|\theta)P(\theta)}{P(S)}$$

Where $S \rightarrow$ training sample with sequence of words $w_1, \dots, w_t < P(w_1), \dots, P(w_k) >$ (where k is vocabulary size)

$\theta \rightarrow$ Set of parameters = for unigram model $P(w_1 | h_1), \dots, P(w_k | h_k) >$ for n-gram

$P(\theta) \rightarrow$ Prior distribution over different possible values of θ

- A point estimate of θ is done by Maximum a Posteriori (MAP) criteria as follows.

$$\theta_{MAP} = \underset{\theta \in \Theta}{\operatorname{argmax}} P(\theta | s) = \underset{\theta \in \Theta}{\operatorname{argmax}} P(S | \theta) P(\theta)$$

Where Θ is space for possible assignments for θ

- The expected value of θ for sample S is given as :

$$\begin{aligned} \theta_B = E[\theta|S] &= \int_{\Theta} \theta P(\theta | S) d\theta \\ &= \frac{\int_{\Theta} \theta P(S | \theta) P(\theta) d\theta}{\int_{\Theta} P(S | \theta) P(\theta) d\theta} \end{aligned}$$

Q.7 Explain large scale language models.

▣ Ans. :

- With the increase of monolingual data scaling of language is required to handle sets of billions or trillions of words.
- In this case exact probability computations and in turn parameter estimation is not feasible.
- So to handle the large data, complete data required for training is divided into partitions.
- Probabilities derived from each partition are stored in separate physical location.
- The language model server handles this data which is distributed over a cluster of independent nodes.
- Clients request statistics from this server during runtime.
- These models facilitate scalability for handling large amount of data, also new data can be added dynamically.
- The disadvantage is slow speed due to networking overheads.

- Another variation can be use of large scale. Distributed language at second pass rescoring stage.
- In this first pass hypothesis is done using smaller language models.
- One more approach is storing large scale, language models in working memory of a single machine.
- The concept of bloom filter is used for this purpose.
- The corpus statistics is stored in memory efficient and randomized data structure called as Bloom filter in quantized manner.

5.6 : Language Model Adaptation

Q.8 What is language mode adaption ? Explain various methods for doing language model adaption.

Ans. : Sometimes of certain domains or languages the training data required to train a model is insufficient.

- In such cases language model adaptation techniques are used to design and tune a language model so that when it is ported to a new domain it will perform well even if very less training data is available.
- One of the commonly used method is model interpolation or mixture languages models, in which the language model is trained with small training data one domain and a generic model is trained using out of domain data.
- These two models are interpolated based on equation.

$$P(w_i | w_{i-1}, w_{i-2}) = \sum_{m=1}^M \lambda_m P(w_i | h_m)$$

- In the second method stated by Seymour and rosenfield the given documents are clustered based on their topics and a language model in built for each cluster. The final model is generated by interpolation of these models. This method is called as topic dependent language model adaptation.
- Dynamic self adaption of a language model is done through trigger models.
- Trigger model consider the word combination likely to co-occur, 10 in turn one word can trigger other for ex. In financial news text the words stock and market are more likely to co-occur.
- For this purpose Lantent Semantic Analysis (LSA) and Probabilistic Latent Semantic Analyses (PLSA) is used.
- The formula for dynamic adaption of a model is

$$P(w_i | h_i, \tilde{h}_i) = \frac{P(w_i | h_i) P(w_i, \tilde{h}_i)}{Z(h_i, \tilde{h}_i)}$$

Where $\tilde{h}_i \rightarrow$ global document history in LSA space up to word i.

- Another approach is unsupervised adaption which is most suitable for speech recognition applications.
- Instead of written text it uses output of speech recognizer.
- One of probability estimation scheme for language model adaption is Maximum Posteriori (MAP) adaption.
- In this the counts are separately collected from the generic out of domain (OD) and in domain adaptation (ID) are combined as follows :

$$P(w|h) = \frac{C_{OD}(w,h) \cdot \epsilon + C_{ID}(w,h)}{C_{OD}(h) \cdot \epsilon + C_{ID}(h)}$$

Where h and w \rightarrow history and predicted words.

5.7 : Types of Language Models

Q.9 Explain class based language model.

▣ Ans. :

- In class based models words are clustered into classes either by automatic means or by linguistic criteria like Part Of Speech (POS) classes.
- To formulate a statistical model it is assumed that words are conditionally independent of other words in a current word class.
- Class based bigram model can be defined as

$$\begin{aligned}
 P(w_i | w_{i-1}) &= \sum_{c_i, c_{i-1}} P(w_i | c_i) P(c_i | c_{i-1}, w_{i-1}) P(c_{i-1} | w_{i-1}) \\
 &= \sum_{c_i, c_{i-1}} P(w_i | c_i) P(c_i | c_i) P(c_{i-1} | w_{i-1})
 \end{aligned}$$

Where c_i is class of word w_i

and c_i is independent of w_{i-1} given c_{i-1}

- Generally a class will contain more than one word, so it can be simplified as

$$P(w_i | w_{i-1}) = P(w_i | c_i) P(c_i | c_{i-1})$$

- By making use of class based models perplexity can be reduced.

Q.10 Explain variable length language models.

▣ Ans. :

- Typically a words are separated by white spaces in the sentences.
- A standard language model considers these vocabulary units to predict the next word based on invariable fixed length history.
- If these units are merged then a type of language model is generated which is called as variable length language model.
- In this approach, initially selected units are merged.
- These merged units which corresponds to frequently observed short phrases are added to language model vocabulary. The example of this can be the words "write" "off" which are different from the word "write-off"
- The probable candidate units are selected based on the mutual information from adjacent words.
- Through greedy iterative algorithm the actual candidate selection is done. The candidates are selected which reduce the perplexity of corpus to maximum extent.
- It is observed that the perplexity is reduced by 10 % compared to word based selection.

Q.11 Explain discriminative language models.

▣ Ans. :

- In the applications like machine translation and speech recognition there is probabtily of generation of sentences which may not match with intended meaning of the input.

- So language model need to do the task of separation of good and bad sentence hypothesis.
- The model should be trained such as the word strings which are significantly different should be assigned distinct probabilities.
- Such models are known as discriminative language models.
- Consider an input x which is a source language string in machine translation application.
- Let y be the set of all possible competing sentences hypothesis.
- Let $GEN(x)$ is a generation function from a language model.
- Let $\phi(x, y)$ is an arbitrary feature function defined jointly over input and each output $y \in y$.
- Based on these functions a global linear modes which selects the best hypotheses is defined as

$$f(x) = \operatorname{argmax}_{y \in GEN(x)} \phi(x, y) \alpha$$

Q.12 What are syntax based language models.

Ans. :

- To understand the need of syntax based language models consider the sentence "Teachers, who sowed the seeds of education in the ancient times, will be remembered for the ages to come".
- In this example the word "Teachers" is responsible for triggering the word "will" but may not be considered by n -gram model as typically the value of n can go to 4 or 5 max.
- These long distance dependencies are addressed by syntax based language models.
- In syntax based language models the syntactic relationships are first modelled and later they are used for estimating better probabilities.
- They consists of statistical parser and a probability model.
- Following are the two approaches for syntax based models.

1. Structured language model it is designed by

s. Chelba and Jelinek

- Joint probability of word sequence and parse s i.e. $p(w, s)$ is computed.
- It is decomposed into product of component probabilities consists of head words from parse structure, POS tags in parse structure and word sequence proper.
- This model when interpolated with trigram model reduced the perplexity by 8% on Wall street journal Continuous Speech Recognition (CSR) and switch board corpora.

2. Almost parsing language model or super AKV model :

- It is designed by Wang and Harper and is based on constraint dependency grammar.
- The sentences are explained with the help of rich tags having syntactic and lexical information of word.
- Joint model called as super ARV language model is defined as follows :

$$P(w_1, \dots, w_n, t_1, \dots, t_n) = \prod_{i=1}^N P(w_i, t_i | w_1, \dots, w_{i-1}, t_1, \dots, t_{i-1})$$

$$= \prod_{i=1}^N P(t_i | w_1, \dots, w_{i-1}) P(w_i | w_1, \dots, w_{i-1}, t_1, \dots, t_{i-1})$$

$$\approx \prod_{i=1}^N P(t_i | w_{i-2}, w_{i-1}, t_{i-1}, t_{i-1}) P(w_i | w_{i-2}, w_{i-1}, t_{i-2}, t_{i-1})$$

Q.13 Explain Max Ent language model

▣ Ans. :

- Maximum likelihood based probabilities estimation for language model face the problem of the constraints as estimation done only of with the help of training data.
- These strong constraints can be relaxed by the use of MaxEnt models.
- Instead of setting probability of given n-gram to its relative frequency in training data the MaxEnt model uses an average observed counts of events in training data.
- The model is formulated as.

$$P(y | x) = \frac{1}{Z(x)} \exp\left(\sum_k \lambda_k f_k(x, y)\right)$$

Where $f(x, y) \rightarrow$ feature function defined on input and predicted variables.

$\lambda \rightarrow$ feature function specific weight

$Z(x) \rightarrow$ normalization factor.

$Z(x)$ can be computed as :

$$Z(x) = \sum_{y \in Y} \exp\left(\sum_k \lambda_k f_k(x, y)\right)$$

Q.14 Explain factored language models.

▣ Ans. :

- Factored language models are based on the observation that if class of the previous word is known, the probability of occurrence of a particular word becomes more accurate.
- If we have the knowledge of class of w_{i-1} word for ex. If the class of w_{i-1} is determiner then we can have good probability estimate for the word w_i i.e. $P(w_i | \text{determiner})$.
- Generalized back off strategy is used for this purpose.
- Words are considered as feature vectors.

$$W = f_1 : k$$

For ex :

Word : stock prices are rising

STEM : stock price be rise

TAG : Nsg N3pl V3pl Vpart

- The statistical model based on this representation is given as

$$P(f_1^{1:k}, f_2^{1:k}, \dots, f_t^{1:k}) \approx \prod_{i=3}^t P(f_i^{1:k} | f_{i-1}^{1:k}, \dots, f_{i-2}^{1:k})$$

- It can be observed from the above equation that each word is not only dependent on single stream of temporally ordered word variables but it is also dependent on additional parallel feature variable.

Q.15 Explain different tree based models.

sub questions can be :

1. Explain hierarchical class based language model.
2. Explain random forest language models.

▣ Ans. :

- The language models which are built using approach using tree structures are known as tree based models.
- Following are some approaches by different researches for construction of tree based models.

Hierarchical class based back off model :

- It is proposed by Zitouni.
- Back off is done keeping general classes at top and specific classes at bottom.
- More specific back off classes are utilized before general once.
- Unlike factored language model in which back off path is combination of probability estimates from different paths and dynamic choice of a path at runtime, the back off path in hierarchical model is fixed.
- This model works better when test data has large number of unseen events.

Model by Wang and Vergyri :

- It is extension of hierarchical class based model.
- In this POS information is used in word clustering.
- Separate hierarchical class trees are defined for various POS categories.

Random Forest Language Models (RFLMs)

- The word histories, present in training data appear as collection of randomly grown decision trees or random forest.
- Root node contain all histories.
- Nodes are associate with set of histories.
- Set of histories are divided into two subsets based on word identity at a particular position in history.
- The split in which log likelihood of training data is maximum is chosen.

Randomness is introduced by following two steps :

- First is random assignment of set of histories from parent node to two sub nodes.
- Set of splits are chosen for log likelihood test randomly.
- When growing process finishes each leaf node can be considered as a cluster of similar word histories which forms equivalence class.
- Decision tree growing procedure is executed multiple times.
- These generated trees are added to random forest.

- For M decision trees RFLM probabilities can be computed as.

$$P_{RF}(w_i | w_{i-n+1}, \dots, w_{i-1}) = \frac{1}{M} \sum_{j=1}^M P_{DT_j}(w_i | \phi_{DT_j}(w_{i-n+1}, \dots, w_{i-1}))$$

Where

$\phi_{DT_j} \rightarrow$ Function which maps the history $w_{i-n+1}, \dots, w_{i-1}$ to a leaf node in j^{th} decision tree.

Q.16 Explain Bayesian topic based language models.

OR Explain Latent Dirichlet Allocation (LDA) model.

▣ **Ans. :**

- Bayesian topic based models are emerged recently in statistical language modeling.
- Blei, Ng and Jordan proposed the first LDA model under Bayesian topic based models.
- LDA model is designed on following assumptions :
 - Each document is considered to be formed of k topics and denoted as Z_1, \dots, Z_k .
 - Topic specific distribution over a particular word is used to generate a word. $k = 1, \dots, k$. probability vector ϕ_k is generated.
 - The prior probabilities of each topic (θ_k). $\theta_1, \theta_2, \dots, \theta_k$ are distributed according to Dirichlet distribution with hyper parameters $\alpha_1, \dots, \alpha_k$ is given as

$$P(\theta_1, \dots, \theta_k) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^k \theta_k^{\alpha_k - 1}$$

- Probability of complete document with sequence W of t words is given as :

$$P(W | \alpha, \phi) = \int P(\theta | \alpha) \left(\prod_{i=1}^t \sum_{Z_i} P(Z_i | \theta) P(w_i | Z_i, \phi) \right) d\theta$$

Q.17 Explain neural network language models.

▣ **Ans. :**

- Neural Network Language Models (NNLMs) are designed based on following approach :
- At first discrete word sequences are mapped to continuous representation.
- Next n-gram probabilities are estimated in continuous space.
- It is assumed that words with similar distribution will have similar continuous representation.
- As shown in Fig. Q.17.1 adjacent layers are fully inter connected.
- For V words vocabulary.
- Input = concatenation of n - 1, V - dimensional binary feature vectors with history of n - 1 words.
- The projection layer i is of fixed d dimensions.
- Projection layer encodes shared continuous representation of words learned during training.
- Hidden layer has fixed no. of J nodes i.e.

$$n_j = \tan h \left(\sum_{k=1}^d w_{jk}^h i_k + b_j^h \right) \forall j, i = 1, \dots, J.$$

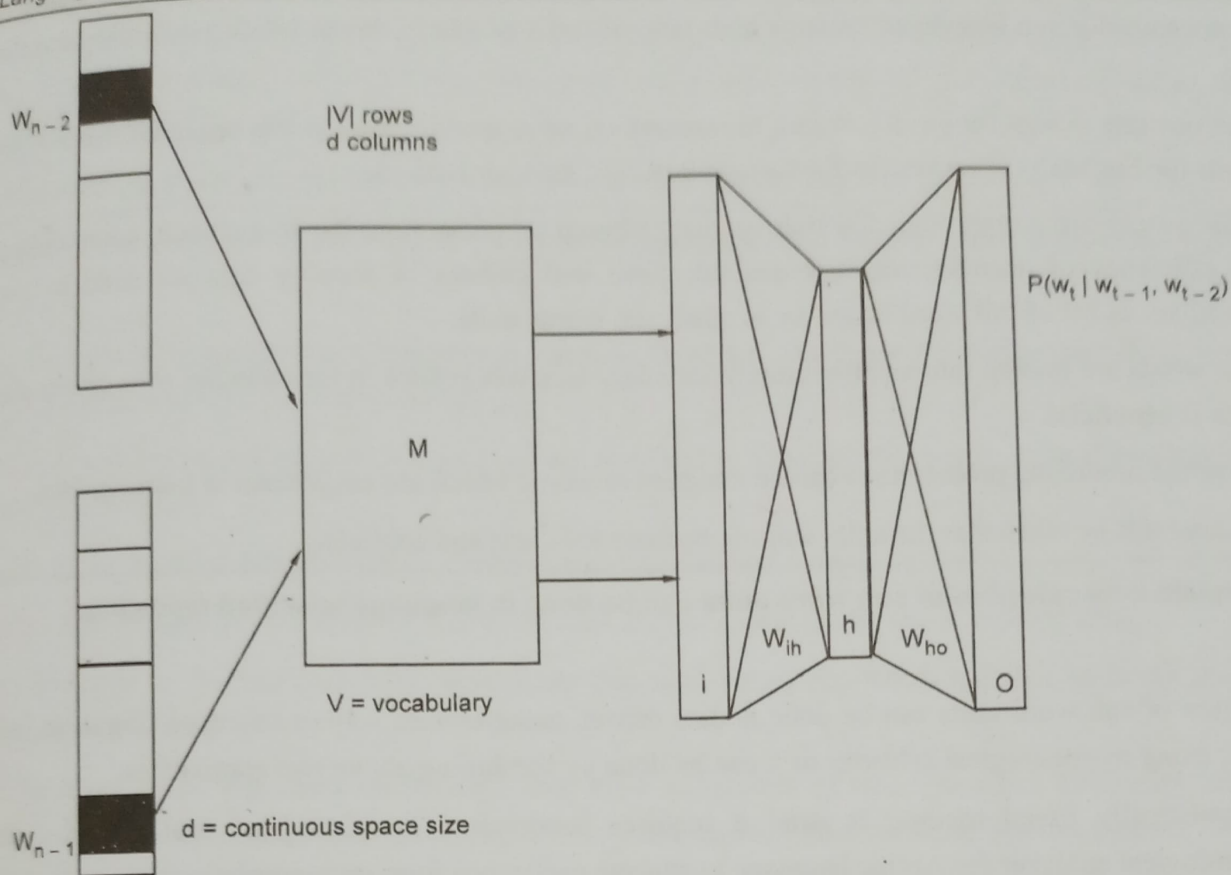


Fig. Q.17.1 : Neural network language model

5.8 : Language Specific Modeling Problems

Q.18 Explain how language modeling is done for morphologically rich languages.

Ans. :

- A morpheme is smallest meaning bearing unit in a language.
- Morphemes are of two types.

1. Free – Occurring on their own
2. Bound – combined with other morpheme.

The process of morphology consists of following :

1. **Compounding** : New word formation from two independent morphemes.
 2. **Derivation** : New work formation by combining free and bound morpheme.
 3. **Inflections** : Obtain a particular grammatical feature by combining free and bound morpheme.
- **Some examples include** : Germanic language with high degree of compounding several morphemes are combined to form a single word in Turkish language, in the languages like Finnish and Arabic root word can have many morphological realization

(for ex : kitaab – iy → my book)

 kitaab – hu → his book

 kitaab – haa → her book ad soon having root word as 'kitaab' with many different realization)

- Morphemically rich languages contains high ratio of types of tokens due to which availability training data is less.
- Also oov rate is high for ex. According to research on news text documents it is observed that oov rate at (N) words for English is 1% where as for Korean it is 25% Turkish 12% etc.
- While processing morphologically rich languages based on constraints due to available computing resources like efficiency of decoder, memory desired speed and amount of training data we need to express the vocabulary as list of full word forms or as small sub words units.
- If the words are broken into smaller units vocabulary size will reduce in turn number of n-grams are reduced which is beneficial.
- Also while modelling probabilities can be assigned to words which are not present in training data.
- But care must be taken that the units will not become too short and confusing.

Q.19 Explain how selection of sub word units can be done in language specified modeling.

▣ Ans. :

- Selection of sub word units can be done is data driven, unsupervised way or based on linguistic information for ex. using morphological analyzer or it can be done by combining above two approaches.
- If linguistically based method is used it requires handcrafted morphological tool for ex. Backwater morphological analyzer for Arabic language to analyze each word form into morphological components.
- If frequency based approach is used it faces the issue of over fitting to the training data and it may hypothesize more morphemes than desired.
- To address this the technique can be used for inclusion of explicit penalty terms for size of morpheme inventory.
- A morpheme inventory M can be derived from corpus C , by maximizing posterior probability as follows :

$$M = \underset{M}{\operatorname{argmax}} P(M | C) = P(C | M) P(M)$$

Q.20 What is the need of modeling a language by morphological categories and how it is done ?

▣ Ans. :

- In any language, while developing a language model it is required to decode interword dependencies as well as dependencies between subword units.
- For this n-gram context need to be increased which can only be done with large training data.
- To address this problem several approaches are stated by researchers.
- One approach is doing the probability assignments based on Statistics over subword components morphological classes.
- One more approach states by Kirchhoff et. al and Vergyri etc. al for Arabic language uses morphological classes like stem root, etc. and words as conditional variables build factored mode.
- This can be further enhanced by addition of additional morphological features to build a neural language model.

- Another approach stated by Sarikaya and Deng on Arabic language uses morphological lexical language model. In this a sentence is annotated with rich parse tree, which contains morphological, syntactic, semantic and other attribute information.
- Random forest language models can be also used for morphological language modelling.

Q.21 Explain how language modeling can be done for the languages without word segmentations.

▣ Ans. :

- There are certain languages like Chinese and Japanese in which words are not segmented at all, instead the document appears as series of character string.
- In Chinese and Japanese language whitespaces are not used to separate the words, instead punctuation signs are used as delimiters.
- One way to model these languages statistically is building a model over characters.
- But it has a disadvantage that inter word relationships are not expressed by the modelling limits.
- So it is desirable to first automatically segment the character string into words and then train the model using these generated words.
- To perform automatic segmentation the algorithm use combination of techniques like dictionary information statistical search, and features like co-occurrence counts non-native letter character positions etc.
- A decoding frame work like viterbi search is used for this purpose.
- Other modelling approaches also include conductional random fields, MaxEnt modelling, discriminative modelling, etc.
- The evaluation of automatic word segmentation done in terms of Precision (P) which indicate percentage of correct cases from all the estimated hypothesized boundaries, R which is percentage of identified boundaries. Out of all possible boundaries
- These two parameters are combined and represented as f measure as :

$$F = 2PR/(P + R)$$

- Performance of 0.96 is achieved through f measure.

Q.22 What are issues faced to device a language modeling system for the languages which does not have a writing system ?

Write a note on language modeling of spoken Vs written languages.

▣ Ans. :

- There are around 6900 languages in the world, out of which many are only spoken languages and does not have electronic presence for ex : Konkani, Arabic.
- We know that for statistical language modelling written data is required.
- So model such languages which are spoken but not written and the languages which are spoken and written but doesn't have orthography i.e. a conventional spelling system is very difficult.
- The difficulties to model the languages with no written script include.
- The language or dialect need to be put into written form manually to generate training data.

- This process is time consuming because of following factors.
 1. Writing standard need to be developed.
 2. To make writing system familiar to the native people, so that they can use it frequently and accurately.
 3. Transcription efforts are required.
- The difficulty to model the languages which lack conventional spelling system is the text which is obtained from web should be normalized which is tedious process.
- Due to above mentioned challenges very less work has been done to model such languages.
- Some of the methods to address this problem which may be used include :
 - Combination of grammar class based approaches with the limited interpreted resources.
 - For the application like developing a dialog system task grammar can be used for predifing the utterances and fine grain word sequence probablities can be modeled from traing a model with the help of less amont of available data.

5.9 : Multilingual and Crosslingual Language Modeling

Q.23 What is the need of multilingual language modeling.

Ans. :

- Typically a standard language model is designed for a particular language.
- In some cases there are some applications like a system which needs to handle the users speaking different languages simultaneously, without knowing which language will be encountered next.
- Another problem which need to be addressed is the system in a specific which should handle different dialects of that language simultaneously.
- One example of language switching is the Hinglish language having the sentences like "I will tell to the teacher ",[®] Lkoky eq'dhy gS", which means in English "I will tell to the teacher that this question is difficult".
- In the above example there is dynamic switching of languages is present.
- To address the problem of dynamic switching between utterances, separate language models can be built from monolingual corper of a language and a system can handle them dynamically based on the highest scoring language model based on output from a first step identification module.
- In the technique stated by fiigen context free grammar having encoded language information in non terminal states and terminal state as n-gram monolingual model can be used ω combine these monolingual models.
- One more way to solve this problem is combining monolingual corpora and train a single model based on this data or interpolate different monolingual language models.
- The switching between the languages or intersentential language switching is a difficult problem due to lack of very less or no training material availability.
- To address this a system is designed by wing et al. which uses four lingual language model with a facility of common back off node as a pause unit.
- Switching between the languages is facilitated after a pause based on the decided probability.

Q.24 What is cross lingual language modeling.

▣ Ans. :

- There are many languages which are low in resourced, i.e. very less or insufficient data is available.
- If we want to work on such languages then to increase probability estimation by extraction of data from other rich resource language.
- One way of using the data from one language to other is to translate it and use it as additional training data.
- Such experiments was done by Kundanpur and kin by translating English text from news domain to generate training data for mandarin language news for speech recognition.
- A unigram is extracted from transited text and is interpolated with trigram baseline model trained on mandarin data . Increased perplexity of 10% is observed.
- Another model is trained by Jensson et. al. for weather reports in Icelandic data using English language data for training purpose. Increased perplexity of 9.20% is observed.
- One approach is to project the source and destination language in common semantic space and construct word translation probabilities by measuring similarity between words from different languages.
- All the above approaches face the problem of dependency on translation accuracy.
- To overcome this problem Tam et al. developed a model which uses Latent Semantic Analysis (LSA) model for source and target language.
- Bilingual Latent Semantic Analysis (BLSA) is done before translation.
- Consider source language as Chinese (ch) and target as English (En).

The probability distribution in English is

$$P_{En}(w) = \sum_k \phi_k^{En}(w) \phi_k^{Ch}$$

$\theta_k \rightarrow$ prior for k^{th} topic

$\phi_k(w) \rightarrow$ Probability assigned to word w by k^{th} latent topic.

- Target language marginal's are then incorporated into target language models as

$$P_{\text{target}}\left(\frac{w}{h}\right) \propto \left(\frac{P_{bLSA}(w)}{P_{\text{base}}(w)}\right)^\beta P_{\text{base}}\left(\frac{w}{h}\right)$$

$P_{bLSA} \rightarrow$ adapted probability

$P_{\text{base}} \rightarrow$ baseline probability

Fill in the Blanks for Mid Term Exam

- Q.1 The _____ algorithms are based on concept of cohesion.
- Q.2 _____ is linking of different textual units by using linguistic devices.
- Q.3 Lexical cohesion is the cohesion which exists in two units which is indicated by _____ between words in those units.
- Q.4 The linguistic expressions like her, she are used to denote an individual is known as _____.