

Universidad Nacional de San Agustín de Arequipa
Escuela Profesional de Ciencia de la Computación

Docentes:

M.Sc. Carlos Eduardo Atencio Torres

Estimación de parámetros (probabilidades de las reglas)

Existen 2 modos

– Usar un corpus de árboles (como el Penn Treebank)

- La probabilidad de una regla se aproxima por la frecuencia relativa de su utilización en el corpus.

La probabilidad de una regla se aproxima por la frecuencia relativa de su utilización en el corpus

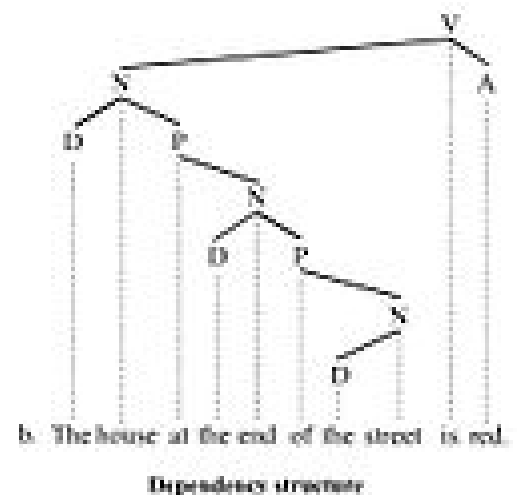
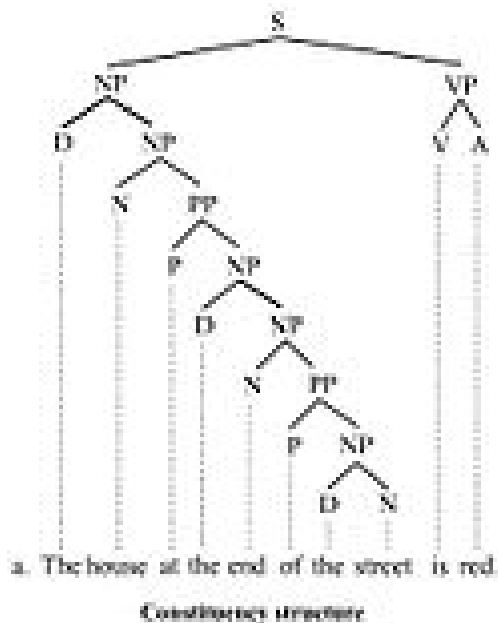
$$P(A \rightarrow \alpha) = C(A \rightarrow \alpha) / \sum \beta C(A \rightarrow \beta) = C(A \rightarrow \alpha) / C(A)$$

$C(A \rightarrow \alpha)$: cantidad de veces que se usa la regla $A \rightarrow \alpha$ en el corpus de árboles

$C(A)$: cantidad de veces que aparece el símbolo A en el corpus

Corpus anotados -Treebanks

- Los treebanks son corpus en los que cada oración está asociada a un árbol sintáctico.
- Se pueden crear: – Directamente, anotando “a mano”. – Con parsing automático y posterior corrección manual.



Penn Treebank

- El Penn Treebank es un corpus anotado ampliamente usado (inglés), mantenido por el LDC (Linguistic Data Consortium).
- Contiene árboles de análisis con información sintáctica y algo de información semántica – una base de datos de árboles lingüísticos.

Categorías gramaticales (tagset)

15 categorías distintas (corpus Brown)

- Incluyen variantes en número para sustantivos
- NN sustantivo singular cat
- NNS sustantivo plural cats

En inglés no hay variación en género (solo pronombres) y ni los adjetivos ni los determinantes varían en número.

- Incluyen variantes en forma, persona y tiempo para verbos
- VB forma base eat
- VBD pasado ate
- VG gerundio eating
- VBN participio eaten
- VBP presente, no 3era persona eat
- VBZ presente, 3era persona eats

Category	Genre (Code)	# of texts	Total Tokens	%
INFORMATIVE	Learned (J)	80	160,000	16.0%
INFORMATIVE	Belles Lettres, Biography, Memoirs, etc (G)	75	150,000	15.0%
INFORMATIVE	Popular Lore (F)	48	96,000	9.6%
INFORMATIVE	Press: Reportage (A)	44	88,000	8.8%
INFORMATIVE	Skills and Hobbies (E)	36	72,000	7.2%
INFORMATIVE	Miscellaneous (H)	30	60,000	6.0%
IMAGINATIVE	General Fiction (K)	29	58,000	5.8%
IMAGINATIVE	Adventure and Western Fiction (N)	29	58,000	5.8%
IMAGINATIVE	Romance and Love Story (P)	29	58,000	5.8%
INFORMATIVE	Press: Editorial (B)	27	54,000	5.4%
IMAGINATIVE	Mystery and Detective Fiction (L)	24	48,000	4.8%
INFORMATIVE	Press: Reviews (theatre, books, music, dance) (C)	17	34,000	3.4%
INFORMATIVE	Religion (D)	17	34,000	3.4%
IMAGINATIVE	Humor (R)	9	18,000	1.8%
IMAGINATIVE	Science Fiction (M)	6	12,000	1.2%
	<i>TOTAL</i>	500	1,000,000	100.0%

Penn Treebank Tagset

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	<i>+, %, &</i>
CD	cardinal number	<i>one, two, three</i>	TO	"to"	<i>to</i>
DT	determiner	<i>a, the</i>	UH	interjection	<i>ah, oops</i>
EX	existential 'there'	<i>there</i>	VB	verb, base form	<i>eat</i>
FW	foreign word	<i>mea culpa</i>	VBD	verb, past tense	<i>ate</i>
IN	preposition/sub-conj	<i>of, in, by</i>	VBG	verb, gerund	<i>eating</i>
JJ	adjective	<i>yellow</i>	VCN	verb, past participle	<i>eaten</i>
JJR	adj., comparative	<i>bigger</i>	VBP	verb, non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>wildest</i>	VBZ	verb, 3sg pres	<i>eats</i>
LS	list item marker	<i>1, 2, One</i>	WDT	wh-determiner	<i>which, that</i>
MD	modal	<i>can, should</i>	WP	wh-pronoun	<i>what, who</i>
NN	noun, sing. or mass	<i>llama</i>	WP\$	possessive wh-	<i>whose</i>
NNS	noun, plural	<i>llamas</i>	WRB	wh-adverb	<i>how, where</i>
NNP	proper noun, singular	<i>IBM</i>	\$	dollar sign	<i>\$</i>
NNPS	proper noun, plural	<i>Carolinas</i>	#	pound sign	<i>#</i>
PDT	predeterminer	<i>all, both</i>	"	left quote	<i>' or "</i>
POS	possessive ending	<i>'s</i>	"	right quote	<i>' or "</i>
PRP	personal pronoun	<i>I, you, he</i>	(left parenthesis	<i>[, (, {, <</i>
PRP\$	possessive pronoun	<i>your, one's</i>)	right parenthesis	<i>],), }, ></i>
RB	adverb	<i>quickly, never</i>	,	comma	<i>,</i>
RBR	adverb, comparative	<i>faster</i>	.	sentence-final punc	<i>. ! ?</i>
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc	<i>: ; ... - -</i>
RP	particle	<i>up, off</i>			

A partir de un corpus se puede definir una gramática, en donde se trata de tomar todos los sub-árboles locales en el conjunto de árboles del corpus.

Usando el Pen podemos inferir una gramática como por ejemplo:

...

(NP

(NP (DT a) (NN round))

(PP (IN of)

(NP

(NP (JJ similar) (NNS increases))

(PP (IN by)

(NP (JJ other) (NNS lenders)))

(PP (IN against)

(NP (NNP Arizona) (JJ real) (NN estate) (NNS loans))))))

...

Reglas "inferidas" de los árboles

NP → NP PP

NP → DT NN

NP → JJ NNS

NP → NP PP PP

$NP \rightarrow NNP \text{ JJ NN NNS}$ (estructuras “chatas”)

$PP \rightarrow IN \text{ NP}$

Hallando las probabilidades

Tenemos una GLC, con reglas como: $NP \rightarrow NP \text{ PP}$

- Extraemos estas reglas del treebank.
- Estimamos las probabilidades de las reglas :
$$p(NP \rightarrow NP \text{ PP}) = \frac{\text{cantidad}(NP \rightarrow NP \text{ PP})}{\text{cantidad}(NP)}$$
- Obtenemos una GLCP.

– Calcularlas sobre un corpus no anotado

- Se comienza con reglas equiprobables.
- Se recalculan las probabilidades según resultados del parsing del paso anterior.
- Se itera hasta converger. (Algoritmo inside-outside)