



UNIVERSIDAD DE BUENOS AIRES  
FACULTAD DE CIENCIAS EXACTAS Y NATURALES  
DEPARTAMENTO DE COMPUTACIÓN

# Nuevas fuentes de información para entrenamiento de etiquetadores gramaticales

Tesis presentada para optar al título de  
Licenciado en Ciencias de la Computación

Fernando Jorge Rodriguez

Director: Dr. José Castaño  
Buenos Aires, 2012

# Índice

<b>1. Introducción</b>	<b>3</b>
1.1. Motivación . . . . .	3
1.2. Trabajo realizado . . . . .	5
1.3. Etiquetado gramatical . . . . .	6
<b>2. Definiciones y marco teórico</b>	<b>7</b>
2.1. Etiquetas . . . . .	7
2.2. Conjuntos de etiquetas . . . . .	11
2.2.1. Especificidad de etiquetas: Treebank vs C5 y Brown . . . . .	13
2.3. Corpus . . . . .	14
2.3.1. Un poco de historia . . . . .	15
2.3.2. Métodos . . . . .	16
2.4. Etiquetadores gramaticales automáticos . . . . .	16
2.4.1. Etiquetadores gramaticales basados en reglas . . . . .	18
2.4.2. Etiquetadores gramaticales estocásticos . . . . .	18
2.4.3. Etiquetadores gramaticales basado en HMM . . . . .	18
2.4.4. Etiquetadores gramaticales de máxima entropía . . . . .	21
2.5. Corpora de entrenamiento y corpora de verificación . . . . .	21
2.6. Evaluación de etiquetadores gramaticales . . . . .	22
2.7. Análisis de error . . . . .	23
2.8. Palabras desconocidas . . . . .	24
2.9. Etiquetador Gramatical TnT . . . . .	25
2.9.1. Modelo teórico . . . . .	25
2.9.2. Suavizado . . . . .	26
2.9.3. Manejo de palabras desconocidas . . . . .	26
2.10. Etiquetador Gramatical Stanford . . . . .	28
2.11. Diccionario COBUILD . . . . .	28
2.11.1. Método de construcción . . . . .	31
2.11.2. Evidencia . . . . .	32
2.11.3. Un corpus . . . . .	32
2.11.4. The Bank of English . . . . .	32
2.11.5. La lista de palabras principales . . . . .	33
2.11.6. Frecuencia . . . . .	33
2.11.7. Ejemplos . . . . .	33
2.11.8. Información gramatical . . . . .	34
2.11.9. Pragmatismo . . . . .	34
2.11.10. Definición del estilo . . . . .	34
2.12. Corpus BNC . . . . .	35
2.13. Corpus WSJ . . . . .	35
<b>3. Desarrollo</b>	<b>35</b>
3.1. Extracción de la información . . . . .	35
3.1.1. Reconocimiento de formas flexionadas . . . . .	36
3.2. Traducción de etiquetas . . . . .	38
3.3. Nuevo Corpus generado . . . . .	40

<b>4. Experimentación</b>	<b>42</b>
4.1. Primer experimento . . . . .	42
4.2. Segundo experimento: Etiquetar el corpus WSJ . . . . .	45
4.2.1. Etiquetar el corpus WSJ con TnT . . . . .	45
4.2.2. Etiquetar el corpus WSJ con Stanford Tagger . . . . .	50
4.3. Tercer experimento: Etiquetar el corpus BNC . . . . .	55
4.3.1. Etiquetar el corpus BNC con TnT . . . . .	55
4.3.2. Etiquetar el corpus BNC con Stanford Tagger . . . . .	60
<b>5. Conclusiones</b>	<b>65</b>

# 1. Introducción

## 1.1. Motivación

El etiquetado o anotado gramatical, también conocido como Part-of-speech tagging, POS tagging o simplemente POST, es el proceso de asignar una etiqueta gramatical a cada una de las palabras de un texto según su categoría léxica. Por ejemplo tomemos la oración siguiente:

(1) *There is no asbestos in our products now.*

El resultado de etiquetarla gramaticalmente es:

(2) *There/EX is/VBZ no/DT asbestos/NN in/IN our/PRP products/NNS now/RB ./.*

donde cada palabra está sucedida por una barra oblicua seguida de la etiqueta gramatical asignada. Se puede apreciar por ejemplo que la palabra *is* fué etiquetada como VBZ (verbo de tiempo presente en tercera persona singular), *products* fué etiquetada como NNS (sustantivo plural), etc. Es decir que a cada palabra se le asignó un código que se corresponde con una función gramatical.

A simple vista el etiquetado gramatical parece una tarea trivial o al menos sencilla, sin embargo no es así. La complejidad de este proceso reside en la ambigüedad gramatical inherente al lenguaje. Por ejemplo, la palabra *premio* puede funcionar como sustantivo:

(1) *Gané un premio*

o como verbo

(2) *Por tu esfuerzo te premio*

En (1), *premio* tendría que recibir la etiqueta gramatical NN (sustantivo común) mientras que en (2) tendría que recibir la etiqueta gramatical VB (verbo). Ahora bien, sería interesante conocer que factor indica cual es la etiqueta correspondiente a una palabra ambigua. En (1) se puede deducir que *premio* es sustantivo porque está precedido por la palabra *un* mientras que en (2) *premio* está precedido por la palabra *te* y a partir de este hecho se puede deducir que en este caso *premio* funciona como verbo. En fin, las palabras circundantes brindan información vital para deducir el sentido gramatical en palabras ambiguas.

Como se menciona más adelante, el etiquetado gramatical juega un papel importante en áreas de la lingüística computacional como síntesis del habla, reconocimiento del habla y recuperación de la información. El etiquetado gramatical es realizado manualmente por lingüistas; especialistas en el lenguaje que se ocupan de determinar una etiqueta gramatical para cada palabra. Desde luego que también es realizado automáticamente por computadoras, mediante programas conocidos como etiquetadores gramaticales. Algunas implementaciones actuales de estos programas están basadas en el aprendizaje; toman un corpus <sup>1</sup>

---

<sup>1</sup>Colección de textos escritos y/o transcripciones del lenguaje oral para cierto idioma

anotado correctamente con el cual se entrenan y luego emplean el conocimiento adquirido para etiquetar el corpus objetivo. En esa primer etapa conocida como entrenamiento, el etiquetador gramatical obtiene, procesa y retiene información sobre cada palabra, su etiqueta asignada y su contexto. Posteriormente el etiquetador determina una etiqueta para cada palabra del corpus objetivo, analizando su ubicación y contexto y utilizando el conocimiento adquirido en la etapa previa.

Uno de los grandes problemas del etiquetado gramatical reside en la falta de corpus anotados para utilizar como corpus de entrenamiento. Los corpus de entrenamiento son etiquetados manualmente por lingüistas especializados. Es un trabajo profundamente meticuloso y tedioso ya que el lingüista debe dar una etiqueta gramatical palabra por palabra en corpus del orden del millón de palabras. Además de la laboriosidad del trabajo, el tiempo empleado para etiquetar un corpus es sumamente extenso y como consecuencia el valor económico es significativo, ya que intervienen grupos de trabajo altamente capacitados durante períodos prolongados. El resultado de este complejo proceso es una tabla de palabras con su correspondiente etiqueta gramatical, como se muestra a continuación:

A	DT
form	NN
of	IN
asbestos	NN
once	RB
used	VCN
to	TO
make	VB
Kent	NNP
cigarette	NN
filters	NNS
has	VBZ
caused	VCN
a	DT
high	JJ
percentage	NN
of	IN
cancer	NN
deaths	NNS
among	IN
a	DT
group	NN
of	IN
workers	NNS
exposed	VCN
to	TO
it	PRP
more	RBR
than	IN
30	CD
years	NNS
ago	RB
,	,
researchers	NNS
reported	VBD
.	.

Ante la importancia que adquieren los corpora etiquetados es inevitable pensar en algún otro tipo de texto que posea información de etiquetas. Por ejemplo algunos diccionarios contienen una palabra, su definición y algunos ejemplos en donde ésta aparece con cada uno de sus sentidos. Es decir que de alguna manera un diccionario contiene por cada palabra uno o más contextos en donde ésta aparece etiquetada. Entonces si tomamos todos los ejemplos de cada palabra de un diccionario y su etiqueta podemos construir un corpus parcialmente anotado. Esta es la idea central de este trabajo.

## 1.2. Trabajo realizado

Como se mencionó en la sección anterior, la idea de este trabajo es suplir la falta de corpus de entrenamiento utilizando la información de etiquetado que

posee un diccionario, generando una nueva fuente de información que servirá para entrenar etiquetadores automáticos. Este trabajo menciona detalladamente la forma de extraer la información relevante sobre etiquetas gramaticales a partir de un diccionario y las decisiones que fueron aplicadas. Esta nueva fuente de información se utiliza para entrenar etiquetadores gramaticales automáticos. Una vez entrenados dichos etiquetadores se emplean para etiquetar un corpus objetivo y se analiza el resultado obtenido. Se realiza el mismo procedimiento, pero ahora combinando la nueva fuente de información generada con un corpus de entrenamiento clásico. Se vuelve a etiquetar un corpus objetivo y se analiza el resultado obtenido. Por último se realizan mediciones sobre el rendimiento de los etiquetadores gramaticales entrenados con esta nueva fuente de información y con los corpora clásicos de entrenamiento y se presentan las conclusiones.

### 1.3. Etiquetado gramatical

Como se mencionó anteriormente, el etiquetado gramatical es el proceso de asignar una etiqueta a cada una de las palabras de un texto según su categoría léxica. Este proceso se realiza en base a la definición de la palabra y la de sus palabras vecinas, es decir, el contexto en que ésta aparece.

Por ejemplo en:

*Does that flight serve dinner*

*dinner* es un sustantivo y por lo tanto recibe la etiqueta para sustantivos NN.

El etiquetado gramatical brinda una gran cantidad de información sobre una palabra y sus vecinas. Por ejemplo, las etiquetas distinguen entre pronombres posesivos (mi, tu, su, etc.) y pronombres personales (Yo, Tú, Él, etc.). Saber si una palabra es un pronombre posesivo o personal nos brinda información sobre las palabras que pueden ocurrir a continuación: los pronombres posesivos generalmente son sucedidos por un sustantivo (como en *Mi comida*) mientras que los personales son sucedidos por un verbo (como en *Yo duermo*).

Utilizando esta deducción podemos aseverar que si una palabra fué etiquetada como pronombre personal, es muy probable que la próxima palabra sea un verbo. Este conocimiento puede ser de útil aplicación en modelos lingüísticos para reconocimiento del habla (voz a texto). Pero esta no es la única información que una etiqueta gramatical puede ofrecer.

Una etiqueta gramatical también nos puede acercar información relacionada con la pronunciación de la palabra. En inglés la palabra *content* puede ser un sustantivo o un adjetivo y su pronunciación varía dependiendo de este hecho. Utilizando estas ideas se pueden producir pronunciaciones más naturales en un sistema de síntesis del habla (texto a voz) o también se puede obtener más exactitud en un sistema de reconocimiento del habla (voz a texto).

Otra aplicación importante del etiquetado gramatical en sistemas de recuperación de la información es el reconocimiento de sustantivos u otro tipo de palabras importantes dentro de un documento, para guardar y utilizar esta información en búsquedas posteriores.

Por último, la asignación automática de etiquetas gramaticales juega un papel importante en algoritmos de desambiguación del sentido de la palabra y en

modelos lingüísticos basados en n-gramas utilizados en sistemas de reconocimiento del habla.

## 2. Definiciones y marco teórico

A continuación se presentan definiciones y teorías que ayudan a comprender el trabajo realizado. Se presenta el concepto de etiqueta gramatical, es decir, una etiqueta que identifica el rol que cumple una palabra dentro de cierto contexto. Se muestran los tipos de etiquetas que han sido utilizados intentando abarcar los distintos significados que pueden tener las palabras. Hasta el día de hoy no se ha llegado a un consenso sobre un conjunto de etiquetas adecuado y se siguen explorando distintas alternativas. Se explica el concepto de etiquetado gramatical, es decir, la tarea de asignar a cada palabra una etiqueta gramatical adecuada según el contexto en donde ésta aparece. Se muestran ejemplos de que esta tarea está muy lejos de ser trivial, introduciendo el concepto de ambigüedad gramatical. Esto ocurre cuando una palabra puede tener muchos significados (y por lo tanto distintas etiquetas gramaticales) dependiendo del contexto en dónde aparece.

Se exhibe la importancia del etiquetado gramatical dentro de distintas áreas como la computación lingüística, reconocimiento y síntesis del habla. Se muestra como se maneja este proceso utilizando programas que lo realizan automáticamente, es decir, etiquetadores gramaticales automáticos. Se explica en profundidad como funcionan estos etiquetadores gramaticales automáticos, mostrando como las implementaciones actuales utilizan un proceso de entrenamiento. Este proceso ocurre a partir de un corpus previamente anotado que el etiquetador automático toma como ejemplo para reproducir el etiquetado.

Se presenta el concepto de corpus y corpus anotados gramaticalmente como conjuntos de información extremadamente valiosos para todas estas tareas. Se muestra la forma de medir, evaluar y comparar el rendimiento de los etiquetadores gramaticales, introduciendo los conceptos de corpus de entrenamiento y corpus de verificación. Se muestran técnicas de análisis de error para el proceso de etiquetación automática. Se exhibe también el manejo de ciertos casos especiales dentro del proceso de etiquetación automático; las palabras desconocidas. Y por último se explican en detalle los etiquetadores automáticos utilizados en el presente trabajo.

### 2.1. Etiquetas

Tradicionalmente la definición de POS o etiqueta gramatical se ha basado en funciones sintácticas y morfológicas, es decir que se agrupan en clases las palabras que funcionan similarmente con respecto a lo que puede ocurrir a su alrededor (sus propiedades de distribución sintáctica) o con respecto a los afijos que poseen (sus propiedades morfológicas). Mientras que las clases de palabras tienen tendencia hacia la coherencia semántica (por ejemplo los sustantivos generalmente describen gente, lugares o cosas y los adjetivos generalmente describen propiedades), este no es necesariamente el caso y en general no se utiliza coherencia semántica como criterio para la definición de POS o etiqueta gramatical.



Las etiquetas gramaticales pueden ser divididas en dos grandes categorías: clases cerradas y clases abiertas. Las clases cerradas son aquellas que tienen miembros relativamente fijos. Por ejemplo, las preposiciones son una clase cerrada porque hay un conjunto cerrado de ellas, es decir que son un grupo de palabras que raramente varía ya que raramente se agregan nuevas preposiciones. En contraste, los sustantivos y los verbos son clases abiertas ya que continuamente se introducen y eliminan nuevos verbos y sustantivos al lenguaje. Es probable que cualquier hablante o corpus tenga una clase abierta de palabras diferente, pero todos los hablantes de un lenguaje y corpora suficientemente grandes, seguramente van a compartir el conjunto de clases de palabras cerradas. Las clases de palabras cerradas también son generalmente palabras funcionales como *de*, *y* o *tu*, que tienden a ser muy cortas, ocurrir frecuentemente y generalmente tienen usos estructurales en gramática.

Hay cuatro clases abiertas principales:

- **Sustantivos** Es el nombre dado a la clase sintáctica que denota personas, lugares o cosas. Pero desde que las clases sintácticas como sustantivos son definidas sintáctica y morfológicamente en vez que semánticamente, algunas palabras para personas, lugares y cosas pueden no ser sustantivos y a la inversa, algunos sustantivos pueden no ser palabras para personas, lugares o cosas. Por lo tanto los sustantivos incluyen términos concretos como *barco* y *silla*, abstracciones como *banda ancha* y *relación*. Se puede definir a una palabra como sustantivo basándose en características como la capacidad de ocurrir con determinantes (una *cabra*, su *banda ancha*), tomar posesivos (los ingresos anuales de *IBM*) y para la mayoría pero no todos los sustantivos, ocurrir en la forma plural (*cabras*, *teléfonos*). Los sustantivos tradicionalmente son agrupados en sustantivos propios y sustantivos comunes.
  - **Sustantivos propios:** Son nombres de personas específicas o entidades y usualmente son escritos en mayúscula.
  - **Sustantivos comunes:** En algunos lenguajes se dividen en sustantivos contables e incontables.
    - **Sustantivos contables:** Son aquellos que permiten establecer su número en unidades. En general esta clase posee forma singular y plural (*silla/s*, *dedo/s*).
    - **Sustantivos incontables:** Se refieren a sustantivos que no se puede determinar su número en unidades (*harina*, *nieve*, *azúcar*).
- **Verbos:** Los verbos son una clase de palabras que incluye a la mayoría de las palabras referidas a acciones y procesos. Tienen ciertas formas morfológicas como tiempo, modo, persona, regularidad, etc. Además, el verbo puede concordar en género, persona y número con algunos de sus argumentos o complementos (a los que normalmente se conoce como sujeto, objeto, etc.). En español concuerda con el sujeto siempre en número y casi siempre en persona (la excepción es el caso del llamado sujeto inclusivo: *Los españoles somos así*).

Algunos ejemplos:

Marisol *canta* una ópera.

La comida *está* caliente.

- **Adjetivos:** Las palabras pertenecientes a esta clase expresan propiedades o cualidades. Por ejemplo *Ese hombre es **alto***. Los adjetivos tienen género y número al igual que los sustantivos. El género y el número de los adjetivos depende del sustantivo al que acompañan. Hay adjetivos que presentan una sola forma para el masculino y para el femenino. Son adjetivos de una sola terminación (verde, especial, amable, grande, etc.). Por el otro lado, los adjetivos de dos terminaciones presentan distintas formas para el masculino y el femenino (feo-fea, pequeño-pequeña, blanco-blanca, etc.) Se clasifican en:

- **Determinativos:** Preceden al sustantivo, lo concretan y lo presentan
  - **Demostrativos:** *Esta* niña
  - **Poseivos:** *Mi* niña
  - **Numerales:** *Tres* niñas
  - **Indefinidos:** *Algunas* niñas
  - **Exclamativos:** ¡*Qué* niña!
  - **Interrogativos:** ¿*Qué* niña?
- **Calificativos:** Califican al sustantivo, es decir, añaden cualidades al sustantivo. Los adjetivos calificativos se dividen en especificativos y explicativos o epítetos.
  - **Adjetivos calificativos especificativos:** Son aquellos que concretan el significado del sustantivo. Suelen aparecer detrás del sustantivo.  
Ej: Quiero una corbata *azul*.
  - **Adjetivos calificativos explicativos o epítetos:** indican cualidades que ya de por sí lleva el sustantivo. Suelen ir delante del sustantivo.  
Ej: *Blanca* nieve, *Verde* hierba.

- **Adverbios:** Los adverbios son otro ejemplo de clase abierta de palabras: se definen como modificadores del verbo, adjetivo o de otro adverbio. Tradicionalmente se dividen en:

- **Adverbios de lugar:** aquí, allí, ahí, allá, acá, arriba, abajo, cerca, lejos, delante, detrás, encima, debajo, enfrente, atrás, alrededor, etc.
- **Adverbios de tiempo absoluto:** pronto, tarde, temprano, todavía, aún, ya, ayer, hoy, mañana, siempre, nunca, jamás, próximamente, prontamente, anoche, enseguida, ahora, mientras.
- **Adverbios de modo:** bien, mal, regular, despacio, deprisa, así, tal, como, aprisa, adrede, peor, mejor, fielmente, estupendamente, fácilmente - todas las que se formen con las terminaciones "mente".
- **Adverbios de cantidad o grado:** muy, poco, muy poco, mucho, bastante, más, menos, algo, demasiado, casi, sólo, solamente, tan, tanto, todo, nada, aproximadamente.

Por otro lado tenemos las clases cerradas de palabras que detallamos a continuación:

- **Preposiciones:** Las preposiciones son enlaces que relacionan los componentes de una oración para brindarles sentido. La unión se lleva a cabo con una o varias palabras. La significación que dan las preposiciones responde a circunstancias de movimiento, lugar, tiempo, modo, causa, posesión, pertenencia, materia y procedencia.  
Algunos ejemplos:

*Me levanté de la cama **a** las ocho de la mañana.*

*Dejé mis cuadernos **sobre** el sillón.*

*Corrí apresurado **hacia** la calle pero no logré divisarte.*

*Lucía se divierte **con** sus muñecas.*

- **Determinantes:** Los determinantes son clases cerradas de palabras que ocurren con sustantivos, generalmente marcando el principio de una frase sustantiva. Un pequeño subtipo de determinantes es el artículo: *a, el*. Otros determinantes incluyen *ese* (como en *el libro ese*).
- **Pronombres:** Los pronombres son formas que generalmente actúan como una clase de atajo para referirse a alguna frase sustantiva, entidad o evento. Se dividen en:
  - **Pronombres personales:** Hacen referencia a personas o entidades (Yo, tú, él, ella, nosotros, ellos, etc.)
  - **Pronombres posesivos:** Son formas de pronombres personales que indican una posesión actual o mas generalmente solo una relacion abstracta entre la persona y algun objeto (mío, tuyo, suyo, mi, nuestro, etc.)
- **Conjunciones:** Las conjunciones son utilizadas para unir dos frases, cláusulas o sentencias. Las conjunciones coordinantes como *y, o* unen dos elementos de igual estado. Las conjunciones subordinativas son utilizadas cuando uno de los elementos es de algún tipo de estado integrado. Por ejemplo *Me molestó **que** no me lo dijeras*.
- **Verbos auxiliares:** Los verbos auxiliares son verbos que proporcionan información gramatical y semántica adicional a un verbo de significado completo. Dichos verbos auxiliares brindan la información gramatical de modo, tiempo, persona y número y las formas no personales. Por ejemplo *¿por qué no **has** llegado a la hora prevista?* o también *La avenida principal de la ciudad **fue** clausurada por obras de refacción*.
- **Numerales:** Los determinantes numerales o simplemente numerales son los que expresan de modo preciso y exacto la cantidad de objetos designados por el sustantivo al que acompañan, delimitan o designan. Limitan el significado general del sustantivo, precisando con exactitud la cantidad de objetos que aquel designa o el lugar de orden que ocupan. Los numerales pueden ser de varias clases. Los más importantes son:

- **Cardinales:** informan una cantidad exacta:  
Quiero *cuatro* libros.
- **Ordinales:** informan del orden de colocación:  
Quiero el *cuarto* libro.
- **Fraccionarios:** informan de particiones de la unidad:  
Quiero la *cuarta* parte.
- **Multiplicativos:** informan de múltiplos:  
Quiero *dobles* ración.

## 2.2. Conjuntos de etiquetas

La sección anterior dió una descripción general de los tipos de clases sintácticas a las que pertenecen las palabras. Esta sección presenta los conjuntos de etiquetas actuales utilizados en la etiquetación gramatical. Es decir, las etiquetas que se corresponden con cada una de estas clases sintácticas. Todavía no existe un consenso sobre el conjunto de etiquetas o tagset más adecuado. Generalmente los conjuntos de etiquetas grandes ofrecen una descripción sintáctica más específica mientras que los conjuntos de etiquetas más pequeños usualmente brindan una información lingüística más acotada. Una de las características clave para decidir que conjunto de etiquetas es el más adecuado justamente depende del nivel de detalle lingüístico que se esté buscando. Otro hecho notable referido a los conjuntos de etiquetas es que los más pequeños generalmente están contenidos en los conjuntos mayores. Ya que las etiquetas más específicas que se encuentran en los conjuntos mayores pueden ser convertidas en etiquetas de menor especificidad con la consecuente pérdida de detalle lingüístico. Por el otro lado, también se pueden convertir las etiquetas pertenecientes a un conjuntos pequeños en etiquetas de mayor especificidad que pertenecen a conjuntos más grandes, ya que generalmente existen etiquetas equivalentes en los conjuntos de mayor tamaño.

**Cuadro 1:** Conjunto de Etiquetas Penn Tree Bank

Etiqueta	Descripción	Ejemplo
CC	Coordinating conjunction	<i>and</i>
CD	Cardinal number	<i>1, third</i>
DT	Determiner	<i>the</i>
EX	Existential	<i>there there is</i>
FW	Foreign word	<i>d'hoivre</i>
IN	Preposition/subordinating conjunction	<i>in, of, like</i>
JJ	Adjective	<i>green</i>
JJR	Adjective, comparative	<i>greener</i>
JJS	Adjective, superlative	<i>greenest</i>
LS	List marker	<i>1)</i>
MD	Modal	<i>could, will</i>
NN	Noun, singular or mass	<i>table</i>
NNS	Noun plural	<i>tables</i>
NNP	Proper noun, singular	<i>John</i>
NNPS	Proper noun, plural	<i>Vikings</i>
PDT	Predeterminer both	<i>the boys</i>

**Cuadro 1:** Conjunto de Etiquetas Penn Tree Bank

Etiqueta	Descripción	Ejemplo
POS	Possessive ending	<i>friend's</i>
PRP	Personal pronoun	<i>I, he, it</i>
PRP\$	Possessive pronoun	<i>my, his</i>
RB	Adverb	<i>however, usually, naturally, here, good</i>
RBR	Adverb, comparative	<i>better</i>
RBS	Adverb, superlative	<i>best</i>
RP	Particle	<i>give up</i>
SYM	Symbol	<i>+, %, &amp;</i>
TO	To	<i>to go, to him</i>
UH	Interjection	<i>uhhuhhuhh</i>
VB	Verb, base form	<i>take</i>
VBD	Verb, past tense	<i>took</i>
VBG	Verb, gerund/present participle	<i>taking</i>
VCN	Verb, past participle	<i>taken</i>
VBP	Verb, sing. present, non-3d	<i>take</i>
VBZ	Verb, 3rd person sing. present	<i>takes</i>
WDT	Wh-determiner	<i>which</i>
WP	Wh-pronoun	<i>who, what</i>
WP\$	Possessive wh-pronoun	<i>whose</i>
WRB	Wh-abverb	<i>where, when</i>
\$	Dollar sign	<i>\$</i>
#	Pound sign	<i>#</i>
"	Left quote	<i>(' or ")</i>
"	Right quote	<i>(' or ")</i>
(	Left parenthesis	<i>( [, (, {, i)</i>
)	Right parenthesis	<i>( ], ), }, i)</i>
,	Comma	<i>,</i>
.	Sentence-final punc	<i>( . ! ?)</i>
:	Mid-sentence punc	<i>( : ; ... -)</i>

Más allá de que no exista aún un consenso sobre que conjunto de etiquetas utilizar, hay un pequeño número de conjuntos de etiquetas o tagsets populares para el idioma inglés, muchos de los cuales evolucionaron a partir del conjunto de etiquetas utilizado para etiquetar el corpus Brown. Este conjunto de etiquetas se conoció como el Brown Corpus Tag-set, un conjunto de 87 etiquetas que se utilizó para etiquetar el corpus Brown: un corpus de 1 millón de palabras construido a partir de ejemplos provenientes de 500 textos de diferentes géneros (diarios, novelas, no ficción, académico, etc.) que fué ensamblado en la Universidad Brown entre 1963 y 1964. Este corpus fué etiquetado gramaticalmente aplicando en primera instancia un etiquetador automático, el programa TAGGIT, y luego corregido manualmente.

Al lado del conjunto de etiquetas Brown se encuentran dos de los conjuntos de etiquetas más utilizados: el conjunto de etiquetas reducido Pen Treebank de 45 etiquetas y el conjunto de etiquetas CLAWS C5 de tamaño medio con 62 etiquetas que fué utilizado para etiquetar el corpus British National Corpus

(BNC).

El conjunto de etiquetas Penn Treebank mostrado anteriormente también fué utilizado para etiquetar el corpus Brown, el corpus Wall Street Journal y el corpus Switchboard entre otros. En realidad, quizás en parte por su pequeño tamaño es uno de los conjuntos de etiquetas más utilizado. A continuación se exhiben algunos ejemplos de oraciones del corpus Brown etiquetadas con el conjunto de etiquetas Penn Treebank. Representaremos una palabra etiquetada mediante la colocación de una barra oblicua seguida de su etiqueta:

1. The/DT grand/JJ jury/NN commented/VBD on/IN a/DT number/NN of/IN other/JJ topics/NNS ./.
2. **There/EX** are/VBP 70/CD children/NNS **there/RB**
3. Although/IN preliminary/JJ findings/NNS were/VBD **reported/VBN** more/RBR than/IN a/DT year/NN ago/IN ./, the/DT latest/JJS results/NNS appear/VBP in/IN today/NN 's/**POS** New/NNP England/NNP Journal/NNP of/IN Medicine/NNP ./.

El primer ejemplo exhibe los determinantes *the* y *a*, los adjetivos *grand* y *other*, los sustantivos comunes *jury*, *number* y *topics* y el verbo en tiempo pasado *commented*. El segundo ejemplo muestra el uso de la etiqueta ET para marcar la construcción existencial *there* y otro uso de *there* que es etiquetado como un adverbio (RB). El tercer ejemplo muestra la segmentación del morfema posesivo 's y un ejemplo de la construcción pasiva 'were reported', en la cual el verbo *reported* está marcado como un pasado participio (VBN) en vez de como un pasado simple (VBD). También es interesante notar que el sustantivo propio *New England* está etiquetado como NNP. Finalmente, se puede observar que como *New England Journal of Medicine* es un sustantivo propio, el etiquetado de Treebank elige marcar cada sustantivo separado como NNP, incluyendo *journal* y *medicine*, que en otro caso serían etiquetados como sustantivos comunes (NN).

### 2.2.1. Especificidad de etiquetas: Treebank vs C5 y Brown

El conjunto de etiquetas Penn Treebank es una selección de 45 etiquetas del conjunto de etiquetas Brown (de 87 etiquetas). Este conjunto reducido deja afuera información que puede ser recuperada desde la identidad del ítem léxico. Por ejemplo los conjuntos de etiquetas Brown y C5 incluyen una etiqueta para cada una de las diferentes formas de los verbos *do*, *be* y *have* (C5 propone la etiqueta VDD para *did* y VDG para *doing*). Estas etiquetas fueron omitidas en el conjunto Treebank.

Ciertas distinciones sintácticas no fueron marcadas en el conjunto de etiquetas Penn Treebank. Por ejemplo, la etiqueta IN es utilizada para preposiciones como para conjunciones subordinadas. El conjunto del Penn Treebank no es suficientemente específico en ciertos casos. Los conjuntos de etiquetas de Brown y C5, por ejemplo, distinguen preposiciones (IN) de conjunciones subordinadas (CS), como en los siguiente ejemplos:

1. **after/CS** spending/VBG a/AT few/AP days/NNS at/IN the/AT Brown/NP Palace/NN Hotel/NN
2. **after/IN** a/AT wedding/NN trip/NN to/IN Corpus/NP Christi/NP ./.

También tienen dos etiquetas para la palabra *to*; en Brown el uso del infinitivo es etiquetado como TO, mientras que las preposiciones son etiquetadas como IN:

1. **to/TO** give/VB priority/NN **to/IN** teacher/NN pay/NN raises/NNS

El conjunto de etiquetas Brown también posee la etiqueta NR para sustantivos adverbiales como *home*, *west*, *Monday* y *tomorrow*. Como Treebank carece de esta etiqueta, hay una política mucho menos consciente para sustantivos adverbiales; *Monday*, *Tuesday* y otros días de la semana son marcados como NNP, *tomorrow*, *west* y *home* son marcados algunas veces como NN y algunas veces como RB. Esto hace al conjunto de etiquetas Treebank menos útil para tareas de alto nivel lingüístico como la detección del tiempo de frases. Sin embargo, el conjunto de etiquetas de Treebank ha sido el más utilizado para la evaluación de algoritmos de etiquetación automática. Esta es la razón por la cual elegimos este conjunto de etiquetas para utilizar en el desarrollo del presente trabajo.

### 2.3. Corpus

Un corpus es una colección de textos escritos y/o transcripciones del lenguaje oral para cierto idioma que generalmente se utiliza para el estudio del lenguaje. La palabra corpus significa cuerpo en latín, su plural es corpora. Habitualmente el tamaño de un corpus es superior al millón de palabras. Para construir un corpus se reúne una cantidad considerable de textos escritos y/o transcripciones orales para luego ser preservado en algún formato (generalmente electrónico).

Los corpora son utilizados por lingüistas para describir naturalmente el lenguaje basados en la evidencia obtenida de sus observaciones. En su trabajo generalmente utilizan operaciones estadísticas sobre los corpora para medir la frecuencia de algún aspecto léxico. Los corpora, grandes cantidades de ocurrencia natural del lenguaje, han ayudado a realizar progresos en diferentes campos del lenguaje como el estudio de fraseología, análisis crítico del discurso, estilos, lingüística forense, traducciones y enseñanza del lenguaje entre otros.

Diferentes tipos de corpora permiten el análisis de distintas clases de discursos para hallar evidencia cuantitativa sobre la existencia de patrones en el lenguaje o para verificar teorías. Los primeros estudios sobre un corpus se enfocaron en palabras; su frecuencia y ocurrencia. Con el desarrollo de la tecnología y de motores de búsqueda más precisos y eficientes, las posibilidades crecieron ampliamente. Hoy en día es posible realizar búsquedas para una palabra perteneciente a cierta clase sintáctica o patrones completos como por ejemplo:

- preposición + sustantivo
- determinante + sustantivo
- una palabra particular + clase de palabra específica sucediéndola.

Cuando corpora escritos y hablados se hicieron disponibles, los lingüistas comenzaron a analizarlos para verificar patrones o diferencias entre el lenguaje hablado y el lenguaje escrito. Parece que aparte de algunas características obvias como salidas en falso y vacilaciones que se producen en el habla, la utilización de un gran número de expresiones deícticas es más frecuente en los discursos orales. Probablemente esto es debido a los signos lingüísticos extra en

los que el lenguaje hablado es más vago. Adicionalmente ciertas características gramaticales manifestadas en el habla deben ser consideradas agramaticales en la escritura.

Otra área importante de estudio lingüístico de corpora es el cambio histórico de los significados de las palabras y la gramática. Y aunque la cantidad de viejos textos disponibles en formato electrónico es mucho más pequeña que la cantidad de textos contemporáneos, el trabajo es factible. En efecto fueron establecidas las diferencias en los aspectos gramaticales concernientes a la voz pasiva.

Por otro lado, en las traducciones es habitual utilizar corpora paralelos que permiten una mejor elección de equivalencias y estructuras gramaticales que podrían reflejar el significado deseado. Estudios adicionales sobre corpora revelaron que los traductores no traducen palabra por palabra sino unidades más grandes (cláusulas o sentencias).

Los estudios de corpora probablemente han tenido una gran influencia en la enseñanza del lenguaje. Primero que nada, han influido en la forma en que se hacen los diccionarios. Segundo los aprendices del lenguaje han sido estudiados para mejorar el conocimiento de los maestros.

Los lingüistas creen que un análisis confiable del lenguaje ocurre mejor en ejemplos recolectados de campo; contextos naturales y con interferencia experimental mínima. Dentro del corpus lingüístico existen visiones divergentes en torno al nivel de las anotaciones. Desde John Sinclair abogando anotaciones mínimas y permitiendo a los textos «hablar por ellos mismos» a otros como el equipo de Survey of English Usage (University College, London) abogando anotaciones como un camino hacia un riguroso entendimiento lingüístico.

### **2.3.1. Un poco de historia**

El punto de inflexión en corpus lingüístico moderno fué la publicación de Henry Kucera y W. Nelson Francis: *Computational Analysis of Present-Day American English* en 1967. Un trabajo basado en el análisis del corpus Brown, una compilación cuidadosamente seleccionada de inglés americano actual totalizando alrededor de 1 millón de palabras obtenidas de una amplia variedad de fuentes. Kucera y Francis sometieron este corpus a una gran variedad de análisis computacional desde el cual compilaron un rico y nutrido corpus combinando elementos de lingüística, enseñanza de lenguaje, psicología, estadística y sociología. Una publicación clave adicional fué «Towards a description of English Usage» de Randolph Quirk (1960) en la que introdujo Survey of English Usage.

Poco después el editor de Boston Houghton-Mifflin se acercó a Kucera para suministrarle el material base de 1 millón de palabras para su nuevo diccionario *American Heritage Dictionary (AHD)*, el primer diccionario que fué compilado utilizando corpus lingüístico. El AHD dió el paso innovador de combinar elementos prescriptivos (como debe utilizarse el lenguaje) con información descriptiva (como se usa actualmente).

Otros editores siguieron el ejemplo. El editor inglés Collins creó y compiló el diccionario *Cobuild* utilizando el corpus *Bank of English*. Fué diseñado para usuarios que están aprendiendo inglés como lengua extranjera.

El corpus Brown también dió lugar a un número de corpora similarmente estructurada: el corpus *LOB* (1960, inglés británico), *Kolhapur* (inglés indio), *Wellington* (inglés de Nueva Zelanda), *Australian Corpus of English* (inglés australiano) y el *Flob corpus* (1990, inglés británico).



Otros corpora representan más lenguajes, variedades y modos: International Corpus of English, el British National Corpus es una colección de 100 millones de palabras provenientes de textos escritos e inglés hablado creado en los 1990s por un consorcio de editores, universidades (Oxford y Lancaster) y la British Library. Para inglés americano contemporáneo, el trabajo se ha centrado en el American National Corpus (más de 400 millones de palabras de inglés americano contemporáneo).

El primer corpus computarizado de lenguaje hablado transcripto fué construido en 1971 por el Montreal French Project, conteniendo 1 millón de palabras que inspiró a Shana Poplack a crear un corpus mucho más grande de Francés hablado.

Al lado de estos corpora de lenguaje vivo se encuentra corpora computarizado que también fué construido a partir de colecciones de textos en lenguajes antiguos. Como ejemplo tenemos la base de datos Andersen-Forbes de la biblia hebrea, desarrollada desde los años 1970, en donde cada cláusula es parseada utilizando grados que representan más de 7 niveles de sintaxis y cada segmento es etiquetado con 7 campos de información. El Quatic Arabic Corpus es un corpus anotado para el lenguaje árabe clásico del corán. Este es un proyecto reciente con múltiples capas de anotación incluyendo segmentación morfológica, etiquetado gramatical y análisis sintáctico utilizando dependencia gramatical.

### 2.3.2. Métodos

Los corpora lingüísticos han generado una cantidad de métodos de investigación intentando trazar un camino desde los datos hacia la teoría. Wallib y Nelson (2001) introdujeron lo que ellos llamaron la perspectiva 3A: anotación, abstracción y análisis.

- **Anotación:** La anotación consiste en la aplicación de un esquema a los textos. Las anotaciones incluyen marcado estructural, etiquetado gramatical, parsing y varias representaciones más.
- **Abstracción:** La abstracción consiste en la traducción (mapeo) de términos del esquema a términos en el modelo teórico. Típicamente incluye búsqueda lingüística directa y también puede incluir aprendizaje por reglas para parsers.
- **Análisis:** El análisis consiste de exploración estadística, manipulación y generalización desde los datos. También podría incluir evaluaciones estadísticas, optimización basada en reglas o métodos de descubrimiento del conocimiento. La mayoría de los corpora de hoy en día está anotado gramaticalmente y aplican algún método para aislar términos que pueden ser interesantes en las palabras circundantes.

## 2.4. Etiquetadores gramaticales automáticos

Como se mencionó anteriormente, el etiquetado gramatical es el proceso de asignar una etiqueta gramatical a cada palabra dentro de un texto. Generalmente las etiquetas gramaticales también son aplicadas a los signos de puntuación, por lo tanto el etiquetado requiere que los signos de puntuación sean separados

de las palabras. Este proceso se realiza previamente o como parte del etiquetado y es conocido como *tokenización*; es el proceso encargado de separar puntos, comas, paréntesis y otros caracteres de las palabras así como también desambiguar el fin de oración (por ejemplo un punto o signo de pregunta) de un signo de puntuación (como en una abreviación por ejemplo *étc.*)

La entrada para un algoritmo de etiquetación automática es una cadena de palabras y un conjunto de etiquetas. La salida es la mejor etiqueta encontrada para cada palabra. Consideremos las siguientes oraciones etiquetadas gramaticalmente:

*Book/VB that/DT flight/NN ./.*

*Does/VBZ that/DT flight/NN serve/VB dinner/NN ?/.*

Asignar una etiqueta gramatical a una palabra no es una tarea trivial incluso en estos sencillos ejemplos. Por ejemplo, la palabra *book* es ambigua. Es decir que tiene más de un uso posible y por lo tanto más de una etiqueta gramatical posible. Puede ser un verbo (como en *book that flight* o *to book the suspect*) o un sustantivo (como en *hand me that book* o *a book of matches*). Análogamente *that* puede ser un determinante (como en *Does that flight serve dinner*) o un complementador (como en *I thought that your flight was earlier*).

El problema del etiquetado gramatical reside en resolver estas ambigüedades, eligiendo la etiqueta adecuada según el contexto. ¿Pero qué magnitud tiene el problema de la ambigüedad de las palabras? Podemos apreciar que la mayoría de las palabras en inglés no son ambiguas, o lo que es lo mismo, tienen una única etiqueta posible. Pero sin embargo muchas de las palabras más comunes del inglés son ambiguas, es decir que las palabras más utilizadas, las que se emplean con mayor frecuencia, pueden tener más de una etiqueta. Por ejemplo *can* puede ser un auxiliar (puede), un sustantivo (lata o contenedor de metal) o un verbo (poner algo en la lata).

Afortunadamente muchas de las palabras ambiguas son fácilmente desambigüables. Esto sucede porque las etiquetas asociadas a una palabra no suelen ocurrir con la misma frecuencia. Por ejemplo *a* puede ser un determinante o la letra *a* (quizás como parte de un acrónimo o una inicial), pero es preciso notar que el sentido de *a* es mucho más frecuente como determinante que como letra. Es decir que es mucho más frecuente encontrar *a* en oraciones como *My father bought a new car* o *There is a hair in my soup* que en oraciones como *Written by A. Kamio* o *The letter a is the first letter of the alphabet*.

Existen distintos métodos computacionales para asignar una etiqueta gramatical a una palabra. La mayoría de los algoritmos de etiquetado automático pertenecen a una de dos clases: etiquetadores basados en reglas o etiquetadores estocásticos.

Los etiquetadores basados en reglas generalmente incluyen una gran cantidad de reglas de desambigüación escritas a mano que especifican, por ejemplo, que una palabra ambigua es un sustantivo antes que un verbo si es seguida por un determinante.

Los etiquetadores estocásticos generalmente resuelven la ambigüedad de etiquetas utilizando un corpus de entrenamiento del cual “aprenden” como etiquetar. Este aprendizaje se realiza extrayendo información sobre la probabilidad de que una palabra dada tenga cierta etiqueta en cierto contexto.

Adicionalmente existe una tercera clase de etiquetadores que es una mezcla

de estos dos: etiquetadores basados en la transformación. Como los etiquetadores basados en reglas, están basados en reglas que determinan cuando una palabra ambigua debe tener cierta etiqueta. Y como los etiquetadores estocásticos tienen un componente de aprendizaje automático; las reglas son inducidas automáticamente a partir de un corpus de entrenamiento previamente etiquetado.

#### **2.4.1. Etiquetadores gramaticales basados en reglas**

Los primeros algoritmos de asignación de etiquetas gramaticales estaban basados en un proceso de dos etapas. En la primer etapa utilizaban un diccionario para asignar a cada palabra una lista de potenciales etiquetas gramaticales. En la segunda etapa utilizaban grandes listas de reglas de desambiguación escritas a mano para reducir la lista de etiquetas hasta llegar a una para cada palabra. De esta manera eliminaban las etiquetas inconsistentes con el contexto.

Las versiones actuales de los etiquetadores gramaticales basados en reglas mantienen los principios originales teniendo en cuenta que los diccionarios y el conjunto de reglas han adquirido un tamaño considerablemente mayor: manejan alrededor de 3800 reglas y un diccionario de etiquetas del orden de las 56.000 entradas para el idioma inglés.

#### **2.4.2. Etiquetadores gramaticales estocásticos**

La inclusión de probabilidades en el proceso de etiquetación gramatical no es una idea nueva. Surge como una consecuencia natural a partir del hecho de que una palabra es empleada con un sentido gramatical mucho más frecuentemente que con otro. Como se mencionó anteriormente, a es mucho más frecuentemente utilizada como determinante que como letra. La inclusión de probabilidades también responde a otro factor importante: la construcción gramatical; cierta etiqueta es precedida frecuentemente por ciertas otra/s. Por ejemplo, como se mencionó anteriormente, los pronombres posesivos generalmente son sucedidos por verbos. Es decir que es más probable encontrar oraciones cuyas palabras estén etiquetadas con PP sucedida por NN que PP sucedida por otra etiqueta.

A continuación vamos a presentar 2 tipos de etiquetadores gramaticales estocásticos: etiquetadores estocásticos basados en el modelo oculto de Markov o simplemente etiquetadores HMM <sup>2</sup> y etiquetadores estocásticos basados en el modelo de máxima entropía.

#### **2.4.3. Etiquetadores gramaticales basado en HMM**

El uso del modelo oculto de Markov para realizar etiquetado gramatical es un caso especial de la inferencia bayesiana, un paradigma que fué conocido a partir del trabajo de Bayes (1763). La inferencia Bayesiana o clasificación Bayesiana fue aplicada exitosamente a problemas del lenguaje a partir de 1950. La clasificación bayesiana puede apreciarse como una tarea para la cual contamos con un conjunto de observaciones y el trabajo consiste en determinar a que conjunto de clases pertenece. En lo que respecta al etiquetado gramatical, se puede utilizar este mismo concepto para tratarlo como una tarea de clasificación de secuencia. En ese caso, la observación será una secuencia de palabras (digamos

---

<sup>2</sup>Por las siglas en inglés de Hidden Markov Model

una oración) para la cual el trabajo consiste en asignar una secuencia de etiquetas gramaticales. Como ejemplo tomemos la oración que aparece a continuación:

*Secretariat is expected to race tomorrow*

En este caso las observaciones son la secuencia de palabras (es decir la oración misma) y nuestro objetivo es asignarles las etiquetas correspondientes. Ya que una palabra puede ser ambigua y tener más de una etiqueta posible, hay una pregunta clave que debemos hacernos: ¿Cuál es la mejor secuencia de etiquetas que corresponden a esta secuencia de palabras? La interpretación bayesiana comienza considerando todas las posibles secuencias de clases –en nuestro caso, todas las posibles secuencias de etiquetas gramaticales. El objetivo aquí es elegir la secuencia de etiquetas que es más probable dada la secuencia de observaciones de  $n$  palabras  $w_1^n$ . En otras palabras, queremos obtener, de todas las secuencias de  $n$  etiquetas  $t_1^n$  la secuencia de etiquetas tal que  $P(t_1^n|w_1^n)$  sea mayor. Se utilizará la notación  $\hat{\phantom{x}}$  para decir “nuestra estimación de la secuencia de etiquetas correcta”.

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} fP(t_1^n|w_1^n) \quad (1)$$

La ecuación anterior se lee así: de todas las secuencias de etiquetas de longitud  $n$ , queremos la secuencia particular  $t_1^n$  que maximiza el lado derecho.

Mientras que esta ecuación nos garantiza obtener la secuencia de etiquetas óptima, todavía no queda del todo claro como utilizarla. Es decir, para una secuencia de etiquetas dada  $t_1^n$  y una secuencia de palabras  $w_1^n$ , no sabemos como computar directamente  $P(t_1^n|w_1^n)$ . Aquí entra en juego la clasificación Bayesiana, ofreciendo una forma de transformar la ecuación en un conjunto de otras probabilidades más sencillas de computar. Las reglas de Bayes reemplazan la probabilidad condicional  $P(x|y)$  por otras tres probabilidades:

$$P(x|y) = \frac{P(y|x)P(x)}{P(y)} \quad (2)$$

Podemos sustituir (2) en (1) para obtener (3):

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} \frac{P(w_1^n|t_1^n)P(t_1^n)}{P(w_1^n)} \quad (3)$$

Convenientemente podemos simplificar (3) eliminando el denominador  $P(w_1^n)$ . Esto sucede ya que estamos eligiendo una de todas las secuencias de etiquetas, computando  $\frac{P(w_1^n|t_1^n)P(t_1^n)}{P(w_1^n)}$  en cada una de ellas. Pero  $P(w_1^n)$  no cambia en ninguna secuencia de etiquetas, entonces estamos preguntando siempre por la misma observación  $w_1^n$ , que tiene la misma probabilidad  $P(w_1^n)$ . Por lo tanto podemos quitar el denominador con la garantía de que el máximo sea el mismo:

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(w_1^n|t_1^n)P(t_1^n) \quad (4)$$

En resumen, la secuencia de etiquetas más probable  $\hat{t}_1^n$  dada alguna palabra  $w_1^n$  puede ser computada tomando el producto de dos probabilidades para cada secuencia de etiquetas y eligiendo la secuencia que lo maximiza.

Desafortunadamente todavía sigue siendo muy difícil computar esta ecuación directamente. Los etiquetadores gramaticales basados en HMM realizan dos suposiciones simplificadoras. La primera es que la probabilidad de aparición de una palabra depende solo de su etiqueta gramatical, es decir que es independiente de las palabras y etiquetas que tiene alrededor. Más técnicamente:

$$P(w_1^n | t_1^n) \approx \prod_{i=1}^n P(w_i | t_i) \quad (5)$$

La segunda suposición es que la probabilidad de aparición de una etiqueta gramatical depende solo de la etiqueta previa (sin tener en cuenta las etiquetas anteriores a la etiquetaa previa), esto es la suposición de bigrama.

$$P(t_1^n) \approx \prod_{i=1}^n P(t_i | t_{i-1}) \quad (6)$$

Utilizando estas suposiciones obtenemos esta nueva ecuación, la cual es utilizada por los etiquetadores gramaticales basados en bigramas para estimar la secuencia de etiquetas gramaticales más probable.

$$\hat{t}_1^n = \underset{t_1^n}{\operatorname{argmax}} P(t_1^n | w_1^n) \approx \underset{t_1^n}{\operatorname{argmax}} \prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1}) \quad (7)$$

La ecuación anterior contiene dos clases de probabilidades, probabilidades de transición de etiquetas y probabilidades de palabras. Tomemos un momento para ver que es lo que representan estas probabilidades.

- **Probabilidades de transición de etiquetas:** Las probabilidades de transición de etiquetas,  $P(t_i | t_{i-1})$ , representan la probabilidad de que ocurra una etiqueta dada la etiqueta previa. Por ejemplo, es muy probable que un determinantes preceda a un adjetivos o a un sustantivo, como *that/DD flight/NN* y *the/DT yellow/JJ hat/NN*. Por lo tanto esperamos que las probabilidades  $P(NN|DT)$  y  $P(JJ|DT)$  sean altas.

Por otro lado, es infrecuente que los adjetivos precedan a los determinantes, entonces la probabilidad  $P(DT|JJ)$  será pequeña. Podemos computar la máxima probabilidad estimada o MLE <sup>3</sup> de una probabilidad de transición de etiquetas  $P(NN|DT)$  etiquetando y contando las etiquetas gramaticales en un corpus. Esto es: de todas las veces que vemos DT, cuántas de esas veces vemos NN después de DT. Lo expresamos más formalmente con el siguiente cociente:

$$P(t_i | t_{i-1}) = \frac{C(t_{i-1}, t_i)}{C(t_i)} \quad (8)$$

Elijamos un corpus específico para examinar, por ejemplo el corpus Brown. Éste es un corpus de 1 millón de palabras de Inglés Americano. El corpus Brown ha sido etiquetado dos veces, la primera en los años sesenta con el conjunto de etiquetas 87-tag y de vuelta en los años noventa con el conjunto de etiquetas Treebank. En el corpus Brown etiquetado con el conjunto

---

<sup>3</sup>Por sus siglas en inglés Maximum Likelihood Estimated

de etiquetas Treebank, la etiqueta DT ocurre 116.454 veces. De esas veces, DT es seguido por NN 56.509 veces. Por lo tanto esta probabilidad de transición se calcula como sigue:

$$P(NN|DT) = \frac{C(DT, NN)}{C(DT)} = \frac{56509}{116454} = .49 \quad (9)$$

Claramente la probabilidad de obtener un sustantivo común después de un determinante es .49 y de hecho alta como sospechábamos.

- **Probabilidades de la palabra:** Por otro lado las probabilidades de la palabra,  $P(w_i|t_i)$ , representan la probabilidad de que dada una etiqueta esta esté asociada con cierta palabra. Por ejemplo si tenemos la etiqueta VBZ (verbo singular de tiempo presente en tercera persona) y quisiéramos adivinar el verbo asociado a esa etiqueta, probablemente elegiríamos el verbo *is*<sup>4</sup>, debido a que el verbo *to be* es muy común en inglés. Podemos computar  $P(is|VBZ)$  de nuevo contando de cuántas veces que vemos VBZ en un corpus cuántas de esas veces VBZ está etiquetando la palabra *is*. Esto es computar el siguiente cociente:

$$P(w_i|t_i) = \frac{C(t_i, w_i)}{C(t_i)} \quad (10)$$

En el corpus Brown etiquetado con Treebank, la etiqueta VBZ ocurre 21.627 veces y VBZ es la etiquetra para *is* 10.073 veces. Entonces:

$$P(is|VBZ) = \frac{C(VBZ, is)}{C(VBZ)} = \frac{10,073}{21,627} = 0,47 \quad (11)$$

Resumiendo, el etiquetado HMM es la tarea de elegir con la mayor probabilidad una secuencia de etiquetas para una secuencia de palabras dada. HMM incluye la suposición de ciertos hechos para simplificar las ecuaciones originales mejorando así la eficiencia de los cálculos.

#### 2.4.4. Etiquetadores gramaticales de máxima entropía

### 2.5. Corpora de entrenamiento y corpora de verificación

Los etiquetadores gramaticales que se basan en modelos estocásticos poseen un proceso de entrenamiento sobre un corpus etiquetado previamente en el cual se generan las probabilidades que se utilizan para tomar decisiones frente a palabras ambiguas.

Dicho corpus de entrenamiento necesita ser cuidadosamente considerado. Si el corpus de entrenamiento es muy específico al dominio, es decir que el corpus de entrenamiento de alguna manera es similar al corpus que se desea etiquetar, las probabilidades van a ser muy ajustadas y no tendrá un buen rendimiento en oraciones de diferentes dominios. Pero si el corpus de entrenamiento es muy general, estas probabilidades no van a llegar a hacer el trabajo suficiente de reflejar el dominio.

---

<sup>4</sup>*is* es el presente en tercera persona del verbo *to be*

Supongamos que estamos intentando etiquetar una oración particular. Si nuestra oración es parte del corpus de entrenamiento, las probabilidades de las etiquetas para esa oración van a ser extraordinariamente precisas y vamos a sobreestimar la precisión de nuestro etiquetador. Se desprende como conclusión que el corpus de entrenamiento no debe ser parcial incluyendo esa oración. Por lo tanto al trabajar con etiquetadores basados en modelos estocásticos, dado un corpus de datos relevante, es una tarea habitual dividir los datos en un corpus de entrenamiento y un corpus de verificación.

Una vez realizada esta división se entrena el etiquetador con el corpus de entrenamiento, se ejecuta el proceso de etiquetación y luego se comparan los resultados con el corpus de verificación.

En general existen dos métodos para entrenar y verificar un etiquetador gramatical. En el primer método, se divide el corpus disponible en tres partes: un corpus de entrenamiento, un corpus de verificación y un corpus de test de desarrollo <sup>5</sup>. Se entrena el etiquetador con el corpus de entrenamiento. Entonces se utiliza el corpus de test de desarrollo para eventualmente afinar o ajustar algunos parámetros y en general decidir cual es el mejor modelo. Una vez que se elige el supuesto mejor modelo, se corre contra el corpus de verificación para ver su rendimiento.

En el segundo método de entrenamiento y verificación, se elige aleatoriamente una división de corpus de entrenamiento y verificación para nuestros datos. Se entrena el etiquetador y luego se calcula el error en el corpus de verificación. A continuación se repite con un corpus de entrenamiento y de verificación diferente seleccionado aleatoriamente. La repetición de este proceso, llamado validación cruzada, generalmente es realizada 10 veces. Luego se promedian esas 10 corridas para obtener un promedio en la proporción del error.

Al comparar modelos es importante utilizar verificaciones estadísticas para determinar si la diferencia entre los modelos es significativa.

## 2.6. Evaluación de etiquetadores gramaticales

Los etiquetadores gramaticales generalmente son evaluados comparando su precisión contra un corpus de verificación <sup>6</sup> etiquetado por humanos. Definimos precisión como el porcentaje de todas las etiquetas en el corpus de verificación donde el etiquetador y el Gold Standard concuerdan. Los algoritmos actuales de etiquetado gramatical tienen una precisión del 96 %-97 % para conjuntos de etiquetas simples como el Penn Treebank. Estas precisiones son para palabras y puntuaciones, la precisión para palabras solas es menor.

Naturalmente uno tiende a preguntarse qué tan bueno es un 97 %. El rendimiento de un proceso de etiquetado puede ser comparado contra un límite inferior y un límite superior. Una manera de establecer un límite superior es ver que tan bien realizan la tarea los humanos.

Marcus, por ejemplo, encontró que los etiquetadores humanos concuerdan en el 96 %-97 % de las etiquetas en el corpus Brown etiquetado con etiquetas Penn Treebank. Esto sugiere que el Gold Standard debe tener un 3 %-4 % de margen de error, y por lo tanto no tiene sentido obtener una precisión del 100 %. Ratnaparkhi mostró que en las palabras donde su etiquetador ha tenido problemas de ambigüedad de etiquetación fueron exactamente las mismas en donde

---

<sup>5</sup>También llamado *devtest*

<sup>6</sup>También llamado *Gold Standard*

los humanos han etiquetado inconsistentemente el corpus de entrenamiento. Dos experimentos realizados por Voutilainen encontraron que cuando a los humanos se les permitió discutir etiquetas, alcanzaron un consenso en el 100 % de las etiquetas.

Por otro lado el límite inferior sugerido por Gale es elegir la etiqueta más probable aplicando el modelo de unigrama para cada palabra ambigua. La etiqueta más probable para cada palabra puede ser computada desde un corpus etiquetado a mano (que puede ser el mismo que el corpus de entrenamiento para el etiquetador que está siendo evaluado).

## 2.7. Análisis de error

Para mejorar el rendimiento de un etiquetador gramatical necesitamos entender donde está funcionando mal. Por eso el análisis de error tiene un papel preponderante. Esta tarea se realiza construyendo una matriz de confusión o tabla de contingencia. Una matriz de confusión es una matriz de  $n \times n$  donde la celda  $(x, y)$  contiene el número de veces que una palabra con correcta etiqueta  $x$  fué etiquetada por el modelo como  $y$ . Por ejemplo, la siguiente tabla muestra una porción de la matriz de confusión para los experimentos de etiquetado con HMM de Franz.

**Cuadro 2:** *Ejemplo de matriz de confusión*

	IN	JJ	NN	NNP	RB	VBD	
IN	-	.2			.7		
JJ	.2	-	3.3	2.1	1.7	.2	2.7
NN		8.7	-				.2
NNP	.2	3.3	4.1	-	.2		
RB	2.2	2.0	.5		-		
VBD		.3	.5			-	4.4
VCN		2.8				2.6	-

Las etiquetas de la fila indican las etiquetas correctas, las etiquetas de las columnas indican las etiquetas asignadas por el etiquetador, y cada celda indica el porcentaje del error de etiquetado general. Por lo tanto 4.4 % del total de errores fueron causados por fallida etiquetacion de VBD como VCN. La matriz anterior y el análisis de error relacionado en Franz, Kupiec y Ratnaparkhi sugieren que algunos de los mayores problemas que encaran los etiquetadores actuales son:

1. NN contra NNP contra JJ: Estas etiquetas son difíciles de distinguir. Es especialmente importante distinguir entre sustantivos propios para extracción de la información y traducción automática.
2. RP contra RB contra IN: Todas estas etiquetas pueden aparecer inmediatamente después del verbo.



3. VBD contra VBN contra JJ: Distinguir estas etiquetas es importante para el *parsing* parcial (los participios son utilizados para encontrar pasivos), y para etiquetar correctamente los bordes de las frases nominales.

El análisis de error es una parte crucial de cualquier aplicación lingüística computacional. Puede ayudar a encontrar *bugs*, encontrar problemas en los datos de entrenamiento y lo más importante, ayuda en el desarrollo de conocimiento y/o algoritmos para utilizar en la solución de problemas.

## 2.8. Palabras desconocidas

Todos los algoritmos de etiquetado gramatical presentados anteriormente requieren un diccionario que liste las posibles etiquetas de cada palabra para que posteriormente el proceso de etiquetado se encargue de identificar la etiqueta correcta. Pero claro, hay un problema: ningún diccionario es capaz de contener todas las palabras. Los sustantivos propios y los acrónimos son creados muy frecuentemente, de hecho ingresan al lenguaje nuevos sustantivos comunes y verbos en una proporción sorprendente. Por lo tanto, para construir un etiquetador completo no podemos utilizar siempre un diccionario para obtener  $P(w_i|t_i)$ . Necesitamos algún método para adivinar la etiqueta de una palabra desconocida.

El algoritmo más básico para manejar palabras desconocidas es suponer que cada palabra desconocida es ambigua entre todas las posibles etiquetas, con igual probabilidad. Entonces el etiquetador debe confiar únicamente en etiquetas contextuales para sugerir la etiqueta adecuada. Un algoritmo ligeramente más complejo está basado en la idea de que la distribución de probabilidad de las etiquetas sobre las palabras desconocidas es muy similar a la distribución de las etiquetas sobre palabras que ocurren solo una vez en un corpus de entrenamiento, una idea sugerida por Baayen y Sproat (1996) y Dermatas y Kokkinakis (1995). Estas palabras que ocurren solo una vez son conocidas como *hapax legomena*.

Por ejemplo, las palabras desconocidas y *hapax legomena* son similares en el hecho de que son más probables de ser sustantivos, seguidas por verbos, pero infrecuentemente suelen ser determinantes o intersecciones. Entonces la probabilidad  $P(w_i|t_i)$  para una palabra desconocida es determinada por el promedio de la distribución sobre todos los conjuntos de palabras de una sola ocurrencia en el corpus de entrenamiento. En resumen, la idea es utilizar “cosas que hemos visto una vez” como un estimador para “cosas que nunca hemos visto”.

De todas maneras, la mayoría de los algoritmos para palabras desconocidas hace uso de una fuente de información mucho más poderosa: la morfología de las palabras. Para el inglés, por ejemplo, palabras terminadas en *s* tienden a ser sustantivos plurales (NNS), palabras terminadas en *ed* tienden a ser pasado participio (VBN), palabras terminadas en *able* tienden a ser adjetivos (JJ), y así. Incluso si nunca vimos una palabra, podemos utilizar hechos sobre su forma morfológica para adivinar su etiqueta. Además la información ortográfica puede ser de mucha ayuda. Por ejemplo, palabras que comienzan con letras mayúsculas generalmente son sustantivos propios (NP). La presencia de un guión es también una característica útil; las palabras con guión en la versión Brown del Treebank son más probables de ser adjetivos (JJ).

¿Cómo son combinadas y utilizadas estas características en los etiquetadores

gramaticales? Un método es entrenar por separado estimadores de probabilidad para cada característica, asumiendo independencia, y multiplicando las probabilidades. Weischedel (1993) construyó un modelo así, basado en cuatro clases específicas. Utilizaron 3 terminaciones inflexionales (*ed*, *s*, *ing*), 32 terminaciones derivacionales (como *ion*, *al*, *ive* y *ly*), 4 valores de mayúscula dependiendo si una palabra es inicio de oración (+/- mayúscula, +/- inicio) y donde una palabra fué guionada. Para cada característica, entrenaron estimadores de máxima verosimilitud de la probabilidad de la característica dada una etiqueta desde un corpus de entrenamiento etiquetado. Entonces combinaron las características para estimar la probabilidad de una palabra desconocida asumiendo independencia y multiplicando:

$$P(w_i|t_i) = p(\text{palabra desconocida}|t_i)p(\text{mayúscula}|t_i)p(\text{final/guión}|t_i) \quad (12)$$

Otro acercamiento basado en HMM, proveniente del trabajo realizado por Samuelsson (1993) y Brants (2000), generaliza el uso de morfología en una manera basada en datos. En este acercamiento, en vez de preseleccionar ciertos sufijos a mano, son consideradas todas las secuencias finales de letras de todas las palabras. Consideran sufijos menores a diez letras, computando para cada sufijo de longitud  $i$  la probabilidad de la etiqueta  $t_i$ :

$$P(t_i|l_{n-i+1}, \dots, l_n) \quad (13)$$

Estas probabilidades son suavizadas utilizando sucesivamente menores y menores sufijos. Esta información de sufijos se mantiene por separado para palabras en mayúscula y minúscula.

En general, la mayoría de los modelos de palabras desconocidas intentan capturar el hecho de que las palabras desconocidas son improbable de ser clases cerradas de palabras. Brants modela este hecho computando solamente las probabilidades de sufijos desde el corpus de entrenamiento para palabras cuya frecuencia en el corpus de entrenamiento es  $\leq 10$ .

## 2.9. Etiketador Gramatical TnT

TnT(Trigrams' n' Tags) es un etiketador gramatical estocástico basado en HMM. Según Brants este etiketador tiene un rendimiento mejor o igual a otros etiketadores actuales de diferentes bases teóricas, incluyendo etiketadores basados en máxima entropía.

### 2.9.1. Modelo teórico

TnT utiliza modelos de Markov de segundo orden para la etiketación gramatical. Técnicamente calcula, dada una secuencia de  $T$  palabras  $w_1, \dots, w_T$

$$\operatorname{argmax}_{t_1, \dots, t_T} \left[ \prod_{i=1}^T P(t_i|t_{i-1}, t_{i-2})P(w_i|t_i) \right] P(t_{T+1}|t_T) \quad (14)$$

para hallar las etiquetas  $t_1, \dots, t_T$ . Las etiquetas adicionales  $t_{-1}, t_0$  y  $t_T$  son delimitadores del principio y el final de la secuencia. Estas etiquetas adicionales mejoran levemente los resultados del etiketado marcando una particularidad de TnT con respecto a otros etiketadores. Las probabilidades son estimadas desde

un corpus etiquetado previamente (el ya mencionado corpus de entrenamiento). Para ello TnT utiliza probabilidades de máxima verosimilitud  $\hat{P}$  obtenidas a partir de la frecuencia relativa y luego aplica una técnica de suavizado

$$\text{Unigramas: } \hat{P}(t_3) = \frac{f(t_3)}{N} \quad (15)$$

$$\text{Bigramas: } \hat{P}(t_3|t_2) = \frac{f(t_2, t_3)}{f(t_2)} \quad (16)$$

$$\text{Trigramas: } \hat{P}(t_3|t_1, t_2) = \frac{f(t_1, t_2, t_3)}{f(t_1, t_2)} \quad (17)$$

$$\text{Léxico: } \hat{P}(w_3|t_3) = \frac{f(w_3, t_3)}{f(t_3)} \quad (18)$$

donde  $t_1, t_2$  y  $t_3$  pertenecen al conjunto de etiquetas y  $w_3$  pertenece al lexicon.  $N$  es el número de *tokens* del corpus de entrenamiento. La probabilidad de máxima verosimilitud se calcula como cero si el denominador o el nominador son cero.

### 2.9.2. Suavizado

TnT aplica una técnica de suavizado sobre las frecuencias contextuales. Esto tiene lugar debido al problema de los datos esparsos en las probabilidades de los trigramas. Es decir, no hay suficientes instancias de cada trigramas para calcular confiablemente su probabilidad asociada. Incluso estableciendo a cero la probabilidad de un trigramas que no aparece en el corpus genera el efecto indeseado de convertir la probabilidad de una secuencia completa en cero. TnT utiliza interpolación lineal de unigramas, bigramas y trigramas para realizar este proceso de suavizado. Es decir que se estima la probabilidad de un trigramas como sigue

$$P(t_3|t_1, t_2) = \lambda_1 \hat{P}(t_3) + \lambda_2 \hat{P}(t_3|t_2) + \lambda_3 \hat{P}(t_3|t_1, t_2) \quad (19)$$

donde  $\hat{P}$  son los estimados de máxima verosimilitud presentados anteriormente y  $\lambda_1, \lambda_2$  y  $\lambda_3$  son los pesos asociados a estos estimadores, tales que  $\lambda_1 + \lambda_2 + \lambda_3 = 1$ . TnT utiliza interpolación lineal con independencia de contexto. Es decir que  $\lambda_1, \lambda_2$  y  $\lambda_3$  tienen el mismo valor para todos los trigramas, o lo que es lo mismo,  $\lambda_1, \lambda_2$  y  $\lambda_3$  son independientes del trigramas que se está calculando. Los valores  $\lambda_1, \lambda_2$  y  $\lambda_3$  son estimados por interpolación de borrado. La idea es que se dará mayor peso a la información de unigramas, bigramas o trigramas más abundante. A continuación se presenta el algoritmo utilizado para realizar esta tarea

### 2.9.3. Manejo de palabras desconocidas

TnT, al igual que muchos otros etiquetadores gramaticales, maneja las palabras desconocidas mediante análisis de sufijos. Los sufijos son fuertes predictores del tipo de palabra. Por ejemplo las palabras terminadas en *able* en el Wall Street Journal parte del Penn Treebank son adjetivos (JJ) en el 98 % de los casos (ej.: *fashionable, variable*) y sustantivos (NN) en el 2 % restante.

La distribución de probabilidades para un sufijo particular es generada a partir de todas las palabras en el corpus de entrenamiento que comparten el

---

**Algoritmo 1** Cálculo de  $\lambda_1, \lambda_2$  y  $\lambda_3 = 0$ 

---

**Establecer**  $\lambda_1 = \lambda_2 = \lambda_3 = 0$   
**por cada** trigramo  $t_1, t_2, t_3$  con  $f(t_1, t_2, t_3) > 0$   
  **según** el máximo de los tres valores siguientes:  
    **caso**  $\frac{f(t_1, t_2, t_3)-1}{f(t_1, t_2)-1}$  : incrementar  $\lambda_1$  en  $f(t_1, t_2, t_3)$   
    **caso**  $\frac{f(t_2, t_3)-1}{f(t_2)-1}$  : incrementar  $\lambda_2$  en  $f(t_1, t_2, t_3)$   
    **caso**  $\frac{f(t_3)-1}{N-1}$  : incrementar  $\lambda_3$  en  $f(t_1, t_2, t_3)$   
  **fin**  
**fin**  
**normalizar**  $\lambda_1, \lambda_2$  y  $\lambda_3$

---

mismo sufijo (de alguna longitud máxima predefinida). El término sufijo se entiende en este contexto como la secuencia final de letras de una palabra, que no coincide necesariamente con el significado lingüístico de sufijo.

La fórmula utilizada para calcular la probabilidad de que una etiqueta pertenezca a cierto sufijo es  $P(t|l_{n-m+1}, \dots, l_n)$ , es decir, la probabilidad de una etiqueta  $t$  dadas las últimas letras  $l_i$  de una palabra de  $n$  letras. TnT aplica una técnica de suavizado utilizando sufijos cada vez más pequeños aplicando un peso  $\theta_i$  a cada uno:

$$P(t|l_{n-m+1}, \dots, l_n) = \frac{\hat{P}(t|l_{n-i+1}, \dots, l_n) + \theta_i P(t|l_{n-i}, \dots, l_n)}{1 + \theta_i} \quad (20)$$

para  $i = m, \dots, 0$ , utilizando el estimador de máxima verosimilitud  $\hat{P}$  para las frecuencias en el lexicon, los pesos  $\theta_i$  y el caso base

$$P(t) = \hat{P}(t) \quad (21)$$

El estimador de máxima verosimilitud para un sufijo de longitud  $i$  es

$$\hat{P}(t|l_{n-i+1}, \dots, l_n) = \frac{f(t, l_{n-i+1}, \dots, l_n)}{f(l_{n-i+1}, \dots, l_n)} \quad (22)$$

TnT utiliza desvío estándar del estimador de máxima verosimilitud para calcular los pesos  $\theta_i$ .

Decisiones de diseño:

1. La primer decisión de diseño que afronta TnT es encontrar un buen valor para  $n$ , la longitud máxima de sufijo utilizada. TnT elige tomar la longitud del mayor sufijo encontrado en el corpus de entrenamiento, con la restricción de que sea menor o igual a 10.
2. Se utiliza independencia de contexto para calcular  $\theta_i$ , la misma idea que se utilizó para calcular  $\lambda_i$ .
3. Se utilizan estimadores distintos para mayúsculas y minúsculas. Es decir, se mantienen dos árboles de sufijos distintos, uno para mayúsculas y otro para minúsculas.
4. La otra decisión relevante es: ¿Qué palabras del lexicon deben ser utilizadas para el manejo de sufijos? Basándose en el hecho de que las palabras

desconocidas son más probablemente infrecuentes, TnT utiliza sufijos de palabras infrecuentes. Por lo tanto, restringe el procedimiento de cálculo de probabilidades de sufijos a palabras con una frecuencia menor o igual a 10.

Adicionalmente, TnT discrimina la información sobre mayúsculas y minúsculas. Esto es debido a que las probabilidades de las etiquetas de palabras con mayúsculas son distintas a las de las palabras con minúsculas. Para llevar esto a cabo se utilizan flags en las probabilidades contextuales. En vez de

$$P(t_3|t_1, t_2) \quad (23)$$

se utiliza

$$P(t_3, c_3|t_1, c_1, t_2, c_2) \quad (24)$$

donde  $c_1$ ,  $c_2$  y  $c_3$  son 1 si la palabra contiene mayúsculas y 0 en otro caso. Esto es equivalente a doblar el conjunto de etiquetas y utilizar etiquetas diferentes según si la palabra aparece en mayúscula o no.

## 2.10. EtiquetadorGramaticalStanford

## 2.11. Diccionario COBUILD

Como se menciona anteriormente, para suplir la falta de corpus de entrenamiento sin caer en la tediosa y costosa tarea de anotar un nuevo corpus, se introduce una fuente de información existente y manualmente anotada. Estamos hablando de un diccionario, que no es ni más ni menos que un conjunto de palabras, donde cada palabra representa una entrada que posee una explicación de su significado, algunas características como su pronunciación y clase gramatical y uno o más ejemplos que muestran su uso. Entonces si pegamos o concatenamos todas esas entradas podemos ver al diccionario como un grupo de palabras con información gramatical para aquellas palabras asociadas a una entrada. Es decir, estas entradas del diccionario:

```
siren
sirens
s*a*!i*%er%e0n
A woman is described as a siren when she is attractive and dangerous to men.
One of the women, another of those sirens, haughtily regarded us as we talked.
countable noun
noun

sirloin
sirloins
s*$e*:l!o!in
A sirloin is a piece of beef which is cut from the lower part of a cows back.
...a sirloin of Scotch beef.
mass noun
noun

sissy
```

sissies

s\*!isi1

A boy is described as a sissy, especially by other boys, if he does not like sport and is  
Youre a lot of cry-babies and sissies... ...Mummys little sissy boy.

countable noun: also vocative

noun

Pueden concatenarse y verse como:

A	is
woman	cut
is	from
described	the
as	lower
a	part
siren NN	of
when	a
she	cows
is	back
attractive	.
and	...
dangerous	a
to	sirloin NN
men	of
.	Scotch
One	beef
of	.
the	Long
women	fibres
,	are
another	picked
of	carefully
those	from
sirens NNS	the
,	sisal NN
haughtily	leaves
regarded	.
us	A
as	boy
we	is
talked	described
.	as
A	a
sirloin NN	sissy NN
is	,
a	especially
piece	by
of	other
beef	boys
which	,

if	.
he	Youre
does	a
not	lot
like	of
sport	cry-babies
and	and
is	sissies NNS
afraid	...
to	...
do	Mummys
things	little
that	sissy NN
are	boy
slightly	.
dangerous	

Esta última información conforma un corpus parcialmente anotado, es decir, un conjunto de oraciones donde alguna/s de las palabras que comprenden cada oración posee/n una etiqueta gramatical. Este corpus parcialmente anotado se utilizará como base para construir un nuevo corpus completamente anotado que servirá como una nueva fuente de información para entrenar etiquetadores gramaticales.

Claramente el primer paso para llevar a cabo esta tarea es elegir un diccionario y extraer la información mencionada anteriormente. El diccionario elegido fué Cobuild. A continuación se detallan las características que lo hicieron distintivo frente a otros diccionarios.

Cobuild es un diccionario basado en la información del corpus Bank of English y el corpus Collins. Su siglas significan: Collins Birmingham University International Language Database. El corpus Collins es una base de datos con alrededor de 2.5 billones de palabras en Inglés. Contiene material escrito de web-sites, diarios, revistas y libros publicados en todo el mundo, y material hablado de radio, TV y conversaciones diarias. A su vez el Bank of English forma parte del corpus Collins. Contiene 650 millones de palabras cuidadosamente seleccionadas para dar un reflejo preciso y balanceado del Inglés que se usa día a día. Gracias a la extensa amplitud del corpus se puede apreciar una gran cantidad de ejemplos de como la gente utiliza realmente el lenguaje. Se puede apreciar el empleo de las palabras, su significado, que palabras ocurren juntas y que tan a menudo. Para decidir que palabras incluir al diccionario Cobuild se ha utilizado información sobre la frecuencia de ocurrencia de las mismas. Por ejemplo, alrededor del 90 % del inglés hablado y escrito está constituido por 3.500 palabras aproximadamente.

The Bank of English contiene un amplio rango de tipos diferentes de lenguaje escrito y hablado proveniente de cientos de fuentes diferentes. Aunque la mayoría de las fuentes son británicas, aproximadamente el 25 % de la información proviene de fuentes de inglés americano y alrededor del 5 % de otras variedades nativas del inglés como Australia y Singapur. Los textos escritos provienen de diarios, revistas, libros de ficción y de no ficción, folletos, informes y cartas. Dos tercios del corpus están confeccionados a partir del lenguaje de los medios: diarios, revistas, radio y televisión. Esta es una categoría significativa

en vista de que millones de personas escuchan y leen el lenguaje presente en los medios. También fueron incluidas publicaciones internacionales, nacionales y locales para capturar un rango general de temas importantes y estilos. Hay otros cientos de libros y revistas de especializadas que abordan temas desde aeróbicos a zoología. Cabe destacar que no fueron incluidos en el corpus libros de texto técnicos, científicos, manuales, etc. El lenguaje hablado informal es representado por grabaciones de conversaciones diarias casuales, reuniones, entrevistas y discusiones. Alrededor de 15 millones de palabras de The Bank of English son transcripciones de lenguaje hablado de esta clase. Luego son seleccionadas para incluir un amplio espectro de temas y situaciones de habla.

El propósito de recolectar toda esta valiosa información en computadoras fué permitirles a los lingüistas (escritores de diccionarios) el acceso a la mayor cantidad de información posible sobre cada una de las palabras que deben definir. Desde luego, los lingüistas son elegidos por su habilidad con el lenguaje, pero ni siquiera el lingüista más hábil puede deducir todos los hechos relevantes sobre las palabras de un lenguaje utilizando solo su intuición. El corpus y el software que se utiliza para analizarlo ayudó al equipo de Cobuild a ordenar la información y ganar valiosa percepción sobre la manera en que se utilizan las palabras: sus significados, sus patrones gramaticales típicos y las maneras en que están relacionadas con otras palabras.

Muchas palabras tienen más de una clase de palabra gramatical asociada y a menudo es de mucha ayuda para los lingüistas mirar solo a una clase de palabra por vez. Para ayudarlos a hacer esto, se ha utilizado un software que muestra las clases de palabras en cada línea del corpus. De esta manera los lingüistas pueden mirar la información completa de la clase de palabra o pueden preguntar solo por verbos, sustantivos, etc. Este tipo de software les permite a los lingüistas tomar decisiones sobre los diferentes sentidos de las palabras, el lenguaje de las definiciones, la selección de ejemplos y la información gramatical dada. El corpus permite realizar esta tarea con confianza y exactitud. Y cuanto más grande es el corpus mayor es la confianza y la exactitud.

El diccionario Cobuild fué concebido teniendo especial atención en los ejemplos expuestos. Como se mencionó anteriormente, el proceso de agregado de palabras al diccionario es muy cuidadoso: cuando un editor quiere agregar una nueva palabra al diccionario, busca en el corpus cada ejemplo que contenga esa palabra. La palabra aparece en una larga lista de oraciones y el editor decide cuál de todos los ejemplos expresa mejor el sentido que está buscando en esa palabra. Todos los ejemplos del diccionario Cobuild muestran patrones gramaticales típicos, vocabulario típico y contextos típicos para cada palabra. En consecuencia, Cobuild presenta una cantidad exhaustiva del vocabulario inglés derivado de observaciones directas del lenguaje.

### **2.11.1. Método de construcción**

En 1987 se publicó el diccionario Cobuild basado en un corpus de 20 millones de palabras. A continuación se construyó un nuevo corpus, el Bank of English con alrededor de 650 millones de palabras. La nueva edición del diccionario Cobuild se basa en este nuevo corpus. La construcción de Cobuild fué un proceso en donde se decidió que palabras y frases presentes en el corpus incluir. Luego se examinó el lenguaje palabra por palabra y frase por frase con el objetivo de dar clara cuenta de cada significado y uso. Entonces para cada entrada se



escribió la definición, se eligieron ejemplos típicos, y se agregó información sobre la pronunciación, la gramática, semántica, pragmatismos y frecuencia.

### **2.11.2. Evidencia**

Un diccionario debe comenzar por la evidencia, los hechos. Los hablantes de un lenguaje conocen mucho sobre éste porque cada día leen y hablan sin esfuerzo durante horas. Sin embargo no son capaces de explicar exactamente que es lo que hacen. La mayoría de las personas no son conscientes de la habilidad que poseen para utilizar un lenguaje; no pueden examinarlo en detalle, simplemente lo utilizan para comunicarse. Aquellos que aprenden a observar el lenguaje cuidadosamente pueden expresar y organizar algunos de los hechos sobre éste basados en la experiencia. De todas maneras hay muchos hechos sobre el lenguaje que no pueden ser descubiertos simplemente pensando y reflexionando sobre él, incluso leyendo y escuchando muy atentamente. Es por eso que Cobuild empleó las computadoras para identificar estos hechos.

### **2.11.3. Un corpus**

El resultado fué que Cobuild estableció un nuevo tipo de evidencia, una colección de textos en inglés llamado corpus ubicado en una computadora de manera que pueda consultarse instantáneamente. Los creadores de Cobuild sabían que necesitaban millones de palabras de inglés, hablado y escrito, americano y británico, formal e informal, sobre hechos y sobre ficción, etc. Esta evidencia reunida durante varios años, permitió encontrar las palabras y expresiones más utilizadas. Cuando una palabra tiene varios significados existe la capacidad de ver cuales son los significados importantes, y que frases se deben incluir. Tomaron como filosofía y fueron conscientes de que todos los detalles de un uso natural de una palabra son esenciales y no pueden ser falsificados. Se dieron cuenta de que debían utilizar ejemplos reales siguiendo la tradición de los grandes lingüistas, en lugar de crearlos.

### **2.11.4. The Bank of English**

Hace varios años que se ha hecho mucho más fácil reunir grandes cantidades de lenguaje hablado y escrito. Los editores de libros, revistas y diarios tomaron consciencia de la gran cantidad de lenguaje que pasaba a través de sus manos y de las muchas buenas razones para conservarlo en formatos electrónicos. De repente apareció un negocio para el lenguaje electrónico entre la gente que quería encontrar o verificar oraciones, particularmente en las noticias, revistas y lenguaje legal. Gradualmente millones de palabras comenzaron a estar disponibles para los estudiosos del lenguaje. Hoy en día el problema no es encontrar el lenguaje sino manejarlo y realizar selecciones sensibles y balanceadas para las tareas analíticas. Diseñando el corpus The Bank of English se balancearon un número de factores(inglés hablado y escrito, americano y británico y otras características: hablantes de comunidades nativas, libros y revistas y más clasificaciones dentro de éstas)

Dentro del componente hablado, el tipo de lenguaje más difícil de recolectar fué como siempre la conversación informal grabada en la vida diaria de la gente común, sin pensar de que su lenguaje está siendo preservado en un corpus.

Cada conversación tiene que ser grabada y transcrita por expertos para luego ser ingresada en una computadora. Esta clase de lenguaje improvisado es de un interés particular para los constructores de diccionarios. El Bank of English cuenta con un total de 15 millones de palabras de este tipo de grabaciones de lenguaje hablado.

#### 2.11.5. La lista de palabras principales

Es mucho más fácil decidir qué palabras y frases incluir y cuales omitir, cuando se tienen cifras exactas sobre una cantidad tan grande de lenguaje. Las computadoras pueden verificar instantáneamente la actividad del lenguaje de miles de hablantes y escritores. Un diccionario (incluso un gran diccionario) es capaz de presentar solo los hechos más importantes del lenguaje y los compiladores necesitan buena evidencia para sus selecciones. Cobuild se especializa en presentar las palabras y frases que son frecuentes en el uso diario. Lejos de ser un registro histórico del lenguaje es más bien una muestra del lenguaje contemporáneo.

#### 2.11.6. Frecuencia

Cobuild brinda información sobre la frecuencia de las palabras principales. Se establecieron 5 bandas de frecuencias. Comenzando con las palabras muy comunes (las de mayor frecuencia), oscila entre un vocabulario básico a uno intermedio hasta cubrir el vocabulario. Las palabras principales sin marca de frecuencia son las menos comunes, sin embargo vale la pena incluirlas en el diccionario. El punto es que el idioma inglés utiliza un número bastante pequeño de palabras para la mayoría de los propósitos pero también tiene disponible un rico y amplio vocabulario. Por ejemplo *be* y *because* pertenecen naturalmente a la banda de mayor frecuencia, por el otro lado, palabras como *barracuda*, *basalt* y *basrelief* no son tan frecuentes. Estas últimas son claramente utilizadas en ocasiones particulares. Cabe aclarar que incluso las palabras infrecuentes incluídas en el diccionario fueron seleccionadas por su utilidad relativa entre miles de palabras posibles.

Entonces Cobuild cuenta con un sistema de frecuencia que marca las palabras principales: una marca significa que la palabra tiene una alta frecuencia y por lo tanto es una palabra común dentro del lenguaje inglés. Dos o más marcas significan que la palabra es parte esencial del vocabulario, cuantas más marcas posee, menos frecuente es.

#### 2.11.7. Ejemplos

Todos los ejemplos fueron seleccionados del corpus The Bank of English. Como se dijo anteriormente, los ejemplos son seleccionados cuidadosamente para mostrar los patrones que aparecen frecuentemente junto a una palabra o frase. El compilador tiene docenas, centenas o miles de ejemplos disponibles y rápidamente escoge los *colocados* <sup>7</sup> y las estructuras típicas en donde la palabra o frase ocurre más a menudo.

Esto significa que los ejemplos cumplen varias funciones. Desde luego ayudan a mostrar el significado de la palabra exhibiendo su uso. Las investigaciones

---

<sup>7</sup>Palabras particulares ubicadas cerca de la palabra principal

incluso sugieren que un gran número de usuarios comienza con los ejemplos antes que con el significado. Las definiciones de Cobuild son bastante claras por sí mismas y los ejemplos muestran el fraseo característico alrededor de la palabra. Como los ejemplos son piezas de texto genuinas y han sido elegidas cuidadosamente en base al uso de la palabra, pueden ser de confianza para exhibir la palabra en un contexto natural.

#### 2.11.8. Información gramatical

Casi cada sentido de cada entrada en el diccionario Cobuild tiene junto a esta una clasificación gramatical, usualmente una clase de palabra y a menudo también una nota estructural. Esta es la información sobre la que se sustenta este trabajo, ya que en base a ella se construirá el nuevo corpus de entrenamiento.

#### 2.11.9. Pragmatismo

Muchos usos de una palabra necesitan más de una frase para explicar apropiadamente su significado. La gente utiliza el lenguaje para realizar muchas cosas: hacer invitaciones, expresar sus sentimientos, enfatizar que es lo que está diciendo, etc. El corpus nos brinda evidencia para tales usos que son difíciles de tomar desde cualquier otra fuente.

El estudio y descripción de las formas en que la gente utiliza el lenguaje para realizar cosas es llamado pragmatismo. Este aspecto del lenguaje es muy importante y fácil de omitir. Esto sucede cuando el lenguaje está dando significado adicional. Cobuild posee mucha información sobre pragmatismo y la expone mediante un símbolo especial en cada entrada. Por ejemplo *and things like that* es definido como una expresión utilizada para ampliar el rango de una lista.

#### 2.11.10. Definición del estilo

La característica más distintiva de Cobuild en su primera versión fué el uso de frases completas en las definiciones. El significado de una palabra fué establecido de la forma en que una persona ordinaria podría explicárselo a otra.

Generalmente los diccionarios ofrecen definiciones breves y tradicionales, mientras que Cobuild expone definiciones realmente amplias y ricas. Si se observan detenidamente las definiciones particulares se puede apreciar que cada palabra es elegida para ilustrar ciertos aspectos del significado. Y en la medida en que es posible, las palabras utilizadas en una definición son más frecuentes que la palabra que está siendo definida.

Las definiciones cortas no pueden decir demasiado. Por ejemplo, el primer sentido de verbo de *mean* podría ser definido como solo *signify*, que es cierto, pero no es todo lo que se puede decir. Cobuild expone esto: *If you want to know...* es decir que ese sentido surge cuando alguien está en la búsqueda de información. La palabra *if* indica que esta es una opción, pero una perfectamente normal, y *you* nos dice que no es una característica de ningún grupo particular de gente (compararlo con *if a policeman arrests you...*). Entonces la definición dice lo que alguien puede querer saber sobre el significado de una *palabra, código, señal o gesto*, indicando que esas son las típicas clases de temas que serán encontradas con este sentido de *mean*. Solo después de toda esta información

viene el equivalente de *signify*: *lo que se refiere a o a que mensaje transmite*. Entonces hay 12 palabras antes de la palabra principal en este sentido, pero cada una de ellas transmite información vital que sería muy difícil de incluir en una definición corta.

## 2.12. Corpus BNC

## 2.13. Corpus WSJ

# 3. Desarrollo

## 3.1. Extracción de la información

El diccionario COBUILD guarda su información en un archivo de texto plano con un formato particular. El primer desafío de este trabajo fué comprender y extraer la información almacenada en ese archivo. A continuación se muestra un pequeño fragmento del mismo para ejemplificar

```
DICTIONARY_ENTRY
ace
aces
*e!*is
If you are or come within an ace of something, you very nearly do or experience it.
He came within an ace of being run over.
phrase: verb inflects
phrase
```

```
DICTIONARY_ENTRY
ace
aces
*e!*is
A person who is ace at something is extremely good at it; an informal use.
...an ace marksman.
classifying adjective
adjective
```

```
DICTIONARY_ENTRY
ace
aces
*e!*is
If you say that something is ace, you mean that you think that it is very good;
an informal use.
Their new records really ace!
qualitative adjective or exclamation
adjective
```

Cada entrada arriba presentada tiene la característica de poseer una cantidad variable de campos y no es posible identificarlos exactamente. Sin embargo, contienen algunos rasgos comunes: la palabra, sus formas, la pronunciación, su definición y uno o más ejemplos donde se indica como se emplea (mediante una etiqueta gramatical). Por ejemplo, en la primer entrada se pueden distinguir

estos campos:

DICTIONARY\_ENTRY

ace → *palabra*

aces → *formas flexionadas*

\*e\*!is → *pronunciación*

If you are or come within an ace of something, you very nearly do or experience it.  
→ *definición*

He came within an ace of being run over. → *ejemplo*

phrase: verb inflects → *etiqueta*

phrase → *etiqueta*

Estas entradas, que conforman el diccionario COBUILD y que constituyen la fuente de información principal sobre la cual se basa este trabajo, fueron cuidadosamente procesadas y refinadas intentando mantener toda la información disponible. El primer desafío de esta etapa consistió en recuperar las entradas con toda la información gramatical disponible; explícita e implícita. Una primer tarea fué reconocer y registrar información relacionada a las formas flexionadas de la palabra (plurales, pasados, etc.), es decir, obtener información gramatical implícita.

### 3.1.1. Reconocimiento de formas flexionadas

En muchas entradas del diccionario COBUILD ocurre la palabra, uno o más ejemplos en donde ésta aparece con cierto sentido (indicado por medio de etiquetas gramaticales) pero dentro de los ejemplos hay apariciones de formas flexionadas. Tomemos la siguiente entrada:

DICTIONARY\_ENTRY

bite → *palabra*

bites, biting, bit, bitten → *formas flexionadas*

b\*a\*!it → *pronunciación*

If an object or surface bites, it grips another object or surface rather than slipping on  
→ *definición*

Let the clutch in slowly until it begins to bite. → *ejemplo*

verb → *etiqueta para la definición*

verb → *etiqueta para el ejemplo*

Aquí arriba se puede observar una entrada del diccionario para la palabra *bite*, que contiene la definición y un ejemplo de esta palabra con sus respectivas etiquetas:

(1) *If an object or surface bites, it grips another object or surface rather than slipping on it or against it.*

(2) *Let the clutch in slowly until it begins to bite.*

En (2) aparece la palabra *bite* en su forma regular con la etiqueta *verb* mientras que en (1) aparece la forma flexionada *bites* con la etiqueta *verb*. En este caso (1) está ofreciendo más información gramatical que la expuesta por medio de

la etiqueta. Reconociendo la forma flexionada (*bites*) podemos adicionarle información extra a la etiqueta *verb*; en vez de guardar la etiqueta de Tree Bank correspondiente a *verb* (VB), en este caso guardaríamos la etiqueta VBZ (verbo de tiempo presente en tercera persona singular) que contiene más información gramatical que VB.

Las entradas de COBUILD exponen las formas derivadas de la palabra que pueden contener los ejemplos. En el ejemplo presentado anteriormente la palabra es *bite* y las formas derivadas de *bite* que muestra la entrada son *bites*, *biting*, *bit* y *bitten*. Con esta información y la etiqueta que fué anotada en COBUILD (*verb*) se puede inferir y generar etiquetas de Tree Bank con información adicional. Como ya se mencionó anteriormente, en este caso la forma *bites* (derivada de la palabra *bite*) que aparece en la definición posee la etiqueta *verb*. La tarea aquí será reconocer que *bites* es un verbo de tiempo presente en tercera persona singular a partir de que *bites* está etiquetada como verbo y de que la palabra de la cual deriva es *bite*. Es decir, inferir el tipo de la forma derivada a partir de la palabra y la etiqueta asignada por COBUILD.

Con el objetivo de identificar las formas derivadas de una palabra se desarrollaron reglas y métodos para su reconocimiento, buscando preservar y aprovechar toda la información que ofrece COBUILD. Entonces, a partir de esta información: la palabra, la forma en que ocurre y la etiqueta asignada se aplican las siguientes reglas para reconocer información adicional a la etiqueta gramatical.

---

**Algoritmo 2** Reconocimiento de formas derivadas

---

Traducir la etiqueta asignada por COBUILD a PenTreeBank

Si la etiqueta obtenida es

**JJ:**

Si la forma termina en *er* o empieza en *more* o *less* aplicar **JJR**

Si la forma termina en *est* o empieza en *most* o *least* aplicar **JJS**

**RB:**

Si la forma termina en *er* o empieza en *more* o *less* aplicar **RBR**

Si la forma termina en *est* o empieza en *most* o *least* aplicar **RBS**

**NN:**

Si la forma termina en *s* aplicar **NNS**

**VB:**

Si la forma termina en *ed* aplicar **VBD—VBN**

Si la forma termina en *ing* aplicar **VBG**

Si la forma termina en *s* aplicar **VBZ**

---

Aplicando algoritmos de extracción y el algoritmo de reconocimiento de formas derivadas explicado anteriormente se obtiene un nuevo corpus parcialmente anotado a partir del diccionario Cobuild. A continuación este corpus será procesado y utilizado como corpus de entrenamiento.

### 3.2. Traducción de etiquetas

Para cada una de sus definiciones, el diccionario COBUILD expone información gramatical expresada mediante etiquetas. Estas etiquetas gramaticales poseen un formato propio. Por ejemplo en la siguiente entrada de COBUILD para la palabra *canary*

```
DICTIONARY_ENTRY
canary
canaries
A canary is a small yellow bird which sings beautifully.
People sometimes keep canaries in cages as pets.
countable noun
noun
```

Se expone la definición (1) y un ejemplo (2), ambos con información gramatical sobre la palabra:

- (1) *A canary is a small yellow bird which sings beautifully.*  
(2) *People sometimes keep canaries in cages as pets.*

Se puede apreciar la etiqueta *noun* asignada por COBUILD para *canary*.

Como la idea de este trabajo es producir un corpus anotado a partir de este diccionario para utilizar como fuente de entrenamiento de etiquetadores gramaticales es necesario que el conjunto de etiquetas empleado sea el mismo que emplea el gold standard para poder medir posteriormente los resultados. Es por eso que se tomó la decisión de traducir estas etiquetas propias de COBUILD en etiquetas de Tree Bank, conjunto con el cual está anotado el gold standard.

A continuación se presenta la tabla de traducción empleada:

**Cuadro 3:** *Tabla de traducción de etiquetas*

Etiqueta COBUILD	Etiqueta PenTreeBank
coordinating conjunction	CC
number	CD
determiner	DT
determiner + countable noun in singular	DT
preposition	IN
subordinating conjunction	IN
preposition, or adverb after verb	IN
preposition after noun	IN
adjective	JJ
classifying adjective	JJ
qualitative adjective	JJ
adjective colour	JJ
ordinal	JJ
adjective after noun	JJ
modal	MD
adverb	RB

**Cuadro 3:** *Tabla de traducción de etiquetas*

Etiqueta COBUILD	Etiqueta PenTreeBank
noun	NN
uncountable noun	NN
noun singular	NN
countable or uncountable noun	NN
countable noun with supporter	NN
uncountable or countable noun	NN
noun singular with determiner	NN
mass noun	NN
uncountable noun with supporter	NN
partitive noun	NN
noun singular with determiner with supporter	NN
countable noun + of	NN
countable noun, or by + noun	NN
countable noun or partitive noun	NN
count or uncountable noun	NN
countable noun or vocative	NN
partitive noun + uncountable noun	NN
noun singular with determiner + of	NN
noun in titles	NN
noun vocative	NN
uncountable noun + of	NN
indefinite pronoun	NN
uncountable noun, or noun singular	NN
countable noun, or in + noun	NN
partitive noun + noun in plural	NN
countable or uncountable noun with supporter	NN
uncountable noun, or noun before noun	NN
uncountable or countable noun with supporter	NN
noun before noun	NN
noun plural with supporter	NNP
noun in names	NNP
proper noun or vocative	NNP
proper noun	NNP
noun plural	NNS
predeterminer	PDT
pronoun	PP
possessive	PPS
adverb with verb	RB
adverb after verb	RB
sentence adverb	RB
adverb + adjective or adverb	RB
adverb + adjective	RB
preposition or adverb	RB
adverb after verb, or classifying adjective	RB
adverb or sentence adverb	RB
adverb with verb, or sentence adverb	RB
exclamation	UH



**Cuadro 3:** Tabla de traducción de etiquetas

Etiqueta COBUILD	Etiqueta PenTreeBank
exclam	UH
verb	VB
verb + object	VB
verb or verb + object	VB
ergative verb	VB
verb + adjunct	VB
verb + object + adjunct	VB
verb + object <i>noun group or reflexive</i>	VB
verb + object or reporting clause	VB
verb + object <i>reflexive</i>	VB
verb + object, or phrasal verb	VB
verb + to-infinitive	VB
ergative verb + adjunct	VB
verb + object + adjunct <i>to</i>	VB
verb + object, or verb + adjunct	VB
verb + object + adjunct <i>with</i>	VB
verb + adjunct <i>with</i>	VB
verb + complement	VB
verb + object, or verb	VB
verb + object + to-infinitive	VB
verb + reporting clause	VB
verb or ergative verb	VB
verb + adjunct <i>from</i>	VB
wh: used as determiner	WDT
wh: used as relative pronoun	WP
wh: used as pronoun	WP
wh: used as adverb	WRB
phrase + noun group	
convention	
combining form	
prefix	
phrasal verb	
other	
phrase	
suffix	
wh	
phrase after noun	
phrase + reporting clause	

### 3.3. Nuevo Corpus generado

A partir del corpus parcialmente anotado obtenido en el proceso de extracción, se completarán las anotaciones automáticamente con un etiquetador gramatical manteniendo las etiquetas gramaticales obtenidas a partir de la información procedente del diccionario COBUILD. Es decir, una vez finalizado el

proceso de extracción de información desde el diccionario, se obtiene un corpus nuevo con las etiquetas gramaticales correspondientes a las palabras definidas en el diccionario. A continuación se exhibe un fragmento del mismo:

```
A
canary NN
is
a
small
yellow
bird
which
sings
beautifully
.
People
sometimes
keep
canaries NNS
in
cages
as
pets
.
```

Este es el resultado de extracción y reconocimiento de formas flexionadas correspondiente a la entrada de COBUILD:

```
DICTIONARY_ENTRY
canary
canaries
A canary is a small yellow bird which sings beautifully.
People sometimes keep canaries in cages as pets.
countable noun
noun
```

Se puede apreciar que se ha reconocido *canaries* como el plural de *canary* (etiqueta NNS) y que se han reconocido y extraído los ejemplos de estas palabras asignando las etiquetas gramaticales traducidas a partir de las etiquetas del diccionario correspondientes a *canary* (countable noun/NN) y *canaries* (noun/NNS).

El próximo paso será el de completar las anotaciones gramaticales para todas las palabras restantes. Este proceso se realiza anotando el corpus plano (sin las etiquetas obtenidas de COBUILD) con el etiquetador gramatical automático TnT. Luego se une este corpus anotado por TnT con el corpus anotado parcialmente procedente de Cobuild, preservando todas las etiquetas del diccionario. El resultado que se muestra a continuación es un nuevo corpus obtenido a partir de Cobuild, con las anotaciones que este provee y completado con anotaciones obtenidas mediante etiquetación automática utilizando TnT.

A	DT
canary	NN
is	VBZ
a	DT
small	JJ
yellow	JJ
bird	NN
which	WDT
sings	VBZ
beautifully	RB
.	.
People	NNS
sometimes	RB
keep	VB
canaries	NNS
in	IN
cages	NNS
as	IN
pets	NNS
.	.

## 4. Experimentación

### 4.1. Primer experimento

El primer experimento consiste en medir (generando una matriz de confusión) la información extraída de COBUILD contra la misma información generada a partir de un etiquetador automático (TnT). Es decir, la información extraída de COBUILD, como se mencionó anteriormente, es la unión de definiciones y ejemplos, con la información gramatical correspondiente a la palabra definida. A continuación se presenta un pequeño extracto:

A	claws
cat NN	that
is	kills
a	smaller
small	animals
furry	such
animal	as
with	mice
a	and
tail	birds
,	.
whiskers	Cats NNS
,	are
and	often
sharp	kept

as	softly
pets	...
.	...
She	domestic
put	animals
out	such
a	as
hand	dogs
and	and
stroked	cats NNS
the	.
cat NN	

Esta es la información extraída de COBUILD para la palabra *cat*; la unión de la definición:

*A cat is a small furry animal with a tail, whiskers, and sharp claws that kills smaller animals such as mice and birds. Cats are often kept as pets.*

y los ejemplos

*She put out a hand and stroked the cat softly...*  
*...domestic animals such as dogs and cats.*

Se puede notar la información gramatical expresada mediante las etiquetas NN y NNS para las palabras *cat* y *cats* respectivamente. La idea de este experimento será comparar estas etiquetas contra las etiquetas asignadas por el etiquetador automático TnT. Entonces se tomará este corpus plano (sin etiquetas), se lo etiquetará utilizando TnT entrenado con el corpus de entrenamiento Wall Street Journal (de ahora en más WSJ) <sup>8</sup> y luego se realizará la comparación.

La matriz de confusión<sup>9</sup> generada a partir de dicha comparación es la siguiente:

**Cuadro 4:** Matriz de confusión para etiquetas extraídas de COBUILD vs generadas por TnT

COBUILD \ TnT	NN	VB	JJ	VCN	RB	VBG	NNP	IN	VBZ	NNS
NN	-	<b>556</b>	<b>1953</b>	52	86	276	-	8	-	-
VB	<b>2616</b>	-	<b>614</b>	-	42	-	77	15	-	5
JJ	<b>1577</b>	96	-	<b>1361</b>	<b>634</b>	<b>555</b>	<b>281</b>	30	-	16
VCN	-	-	-	-	-	-	-	-	-	-
RB	219	23	<b>408</b>	10	-	9	34	249	-	11
VBG	-	-	-	-	-	-	-	-	-	-
NNP	-	-	-	-	-	-	-	-	-	-
IN	-	-	-	-	-	-	-	-	-	-
VBZ	-	-	-	-	-	-	-	-	-	-
NNS	83	1	17	-	1	2	104	3	192	-

<sup>8</sup>Wall Street Journal es un corpus anotado, parte del Penn Treebank

<sup>9</sup>Las matrices de confusión presentadas de aquí en adelante contienen las primeras 10 etiquetas de mayor error

Porcentaje de aciertos: 99,16 %  
Cantidad de errores: 13082

Se puede apreciar un alto porcentaje de aciertos entre las etiquetas extraídas de COBUILD (99,16 %) y las etiquetas asignadas por TnT. Este porcentaje indica que la información de etiquetas extraídas de COBUILD es consistente con las producidas por TnT. La mayoría de los errores se da en etiquetas VB, NN y JJ de COBUILD cuando son etiquetadas como NN, JJ y NN por TnT respectivamente.

## 4.2. Segundo experimento: Etiquetar el corpus WSJ

El segundo experimento realizado tiene como objetivo evaluar la nueva fuente de información obtenida (NFI) como corpus de entrenamiento. Para esto se entrenarán 2 etiquetadores gramaticales (Stanford Tagger y TnT) y se etiquetará con ellos el corpus Wall Street Journal (WSJ). Posteriormente se realizarán mediciones de desempeño pertinentes.

### 4.2.1. Etiquetar el corpus WSJ con TnT

La primer evaluación de este segundo experimento consiste en entrenar el etiquetador gramatical TnT con WSJ como corpus de entrenamiento y con WSJ + NFI. Luego se procede a etiquetar el WSJ plano (sin etiquetas gramaticales) con estos dos modelos. Por último se contruye la matriz de confusión:

**Cuadro 5:** *WSJ original contra WSJ etiquetado con TnT (entrenado con WSJ)*

TnT(WSJ) WSJ	JJ	NNP	VBN	RB	IN	RP	NNPS	VBD	NN	VB
<b>NN</b>	<b>2592</b>	<b>1649</b>	28	79	18	7	-	39	-	316
<b>VBD</b>	70	5	<b>1642</b>	-	-	-	-	-	42	41
<b>IN</b>	108	30	-	<b>1476</b>	-	<b>1376</b>	-	-	4	2
<b>RB</b>	726	45	4	-	<b>1459</b>	857	-	1	200	37
<b>VBN</b>	<b>1262</b>	22	-	-	-	-	-	<b>1146</b>	36	53
<b>NNP</b>	484	-	4	46	23	1	<b>1236</b>	3	378	7
<b>JJ</b>	-	699	<b>916</b>	601	69	23	1	42	903	70
<b>VBG</b>	424	20	-	-	-	-	-	-	884	-
<b>VBP</b>	26	10	33	11	10	1	1	51	349	872
<b>WDT</b>	-	-	-	-	722	-	-	-	-	-

Aciertos: 1.226.484 (97,10 %)

Errores: 36.636 (2,90 %)

**Cuadro 6:** *WSJ original contra WSJ etiquetado con TnT (entrenado con WSJ + NFI)*

TnT(WSJ+NFI) WSJ	JJ	NNP	VBN	IN	RP	NNPS	RB	VBD	VB	NN
<b>NN</b>	<b>2794</b>	<b>2081</b>	46	16	7	-	87	43	315	-
<b>VBD</b>	73	10	<b>1716</b>	-	-	-	-	-	37	22
<b>RB</b>	799	62	4	<b>1664</b>	<b>1128</b>	-	-	1	50	204
<b>IN</b>	107	51	-	-	<b>1559</b>	-	<b>1222</b>	-	3	4
<b>VBN</b>	<b>1369</b>	27	-	-	-	-	-	<b>1219</b>	42	33
<b>NNP</b>	358	-	4	19	1	<b>1257</b>	44	3	9	263
<b>JJ</b>	-	1013	893	87	28	1	634	47	74	864
<b>VBP</b>	30	12	33	10	1	1	8	56	885	379
<b>VBG</b>	518	27	-	-	-	-	-	-	-	816
<b>VBZ</b>	-	11	-	-	-	2	-	-	-	1

Aciertos: 1.226.967 (97,14 %)

Errores: 36.153 (2,86 %)

Se puede observar que el rendimiento del etiquetador TnT entrenado con WSJ+NFI es un poco mejor (97,14 %) que el rendimiento de TnT entrenado con WSJ (97,1 %). La mayoría de los errores para TnT entrenado con WSJ se da en etiquetas NN del gold standard cuando son etiquetadas como JJ y NNP por TnT. Para TnT entrenado con WSJ + NFI la mayoría de los errores se da en las mismas etiquetas, pero con cantidad de errores mayor, sobre todo para NN etiquetado como NNP.

La segunda evaluación de este experimento consiste en entrenar TnT con la mitad de WSJ y con la mitad de WSJ + NFI. Posteriormente con estos dos modelos se etiqueta la mitad restante de WSJ y se construye la matriz de confusión. Se realiza la misma operación para cada mitad:

**Cuadro 7:** 1 mitad WSJ original contra 1 mitad WSJ etiquetado con TnT (entrenado con 2 mitad de WSJ)

TnT(2WSJ) WSJ	JJ	NNP	VBN	NN	VBD	IN	RB	RP	VB	NNPS
<b>NN</b>	<b>1959</b>	<b>1154</b>	26	-	24	5	60	2	269	2
<b>VBD</b>	76	12	<b>1129</b>	19	-	-	1	-	29	-
<b>JJ</b>	-	545	<b>801</b>	<b>1039</b>	52	19	313	9	62	1
<b>VBN</b>	<b>617</b>	23	-	24	<b>819</b>	-	-	-	36	-
<b>RB</b>	432	25	3	91	2	<b>808</b>	-	318	19	-
<b>IN</b>	71	24	1	3	-	-	<b>634</b>	<b>615</b>	1	-
<b>VBP</b>	26	19	19	285	33	6	4	-	613	1
<b>NNP</b>	419	-	8	534	11	19	43	-	20	600
<b>VBG</b>	276	22	-	577	-	-	-	-	-	-
<b>NNPS</b>	26	549	-	-	-	-	-	-	-	-

Aciertos: 607.876 (96,25 %)

Errores: 23.695 (3,75 %)

**Cuadro 8:** 1 mitad WSJ original contra 1 mitad WSJ etiquetado con TnT (entrenado con 2 mitad de WSJ + NFI)

TnT(2WSJ+NFI) WSJ	JJ	NNP	VBN	IN	NN	RP	VBD	NNPS	VBG	VB
<b>NN</b>	<b>1759</b>	<b>1287</b>	29	6	-	3	27	1	556	213
<b>VBD</b>	54	17	<b>1039</b>	-	10	-	-	-	-	17
<b>RB</b>	434	30	2	<b>872</b>	81	511	1	-	-	23
<b>JJ</b>	-	<b>689</b>	<b>612</b>	37	<b>838</b>	6	29	-	219	45
<b>IN</b>	65	33	-	-	3	<b>749</b>	-	-	2	1
<b>VBN</b>	<b>654</b>	24	-	-	16	-	<b>708</b>	-	-	21
<b>NNP</b>	334	-	6	15	356	-	4	558	23	18
<b>VBP</b>	14	19	20	6	248	-	28	1	-	534

**Cuadro 8:** 1 mitad WSJ original contra 1 mitad WSJ etiquetado con TnT (entrenado con 2 mitad de WSJ + NFI)

$\begin{matrix} \text{TnT(2WSJ+NFI)} \\ \text{WSJ} \end{matrix}$	JJ	NNP	VBN	IN	NN	RP	VBD	NNPS	VBG	VB
<b>NNPS</b>	20	510	-	-	1	-	-	-	-	-
<b>VBG</b>	322	22	-	-	460	-	-	-	-	-

Aciertos: 609.255 (96,47 %)

Errores: 22.316 (3,53 %)

**Cuadro 9:** 2 mitad WSJ original contra 2 mitad WSJ etiquetado con TnT (entrenado con 1 mitad de WSJ)

$\begin{matrix} \text{TnT(1WSJ)} \\ \text{WSJ} \end{matrix}$	JJ	VBN	NNP	NN	RB	VBD	NNPS	IN	RP	VB
<b>NN</b>	<b>1826</b>	35	<b>1089</b>	-	51	27	-	11	6	256
<b>VBD</b>	73	<b>1097</b>	12	37	-	-	-	-	-	19
<b>JJ</b>	-	609	559	<b>1085</b>	360	69	2	44	18	78
<b>IN</b>	44	-	19	6	<b>881</b>	-	-	-	<b>693</b>	4
<b>VBN</b>	<b>838</b>	-	30	22	-	<b>859</b>	-	-	-	33
<b>NNP</b>	457	17	-	458	40	6	<b>855</b>	20	2	18
<b>RB</b>	384	3	45	163	-	-	-	<b>741</b>	517	19
<b>VBP</b>	35	17	14	187	8	31	-	7	1	560
<b>VBG</b>	294	-	21	552	1	-	-	-	-	1
<b>VB</b>	74	25	49	405	9	23	-	7	-	-

Aciertos: 607.593 (96,21 %)

Errores: 23.956 (3,79 %)

**Cuadro 10:** 2 mitad WSJ original contra 2 mitad WSJ etiquetado con TnT (entrenado con 1 mitad de WSJ)

$\begin{matrix} \text{TnT(1WSJ+NFI)} \\ \text{WSJ} \end{matrix}$	JJ	NNP	VBN	RP	IN	NN	NNPS	VBD	RB	VBG
<b>NN</b>	<b>1831</b>	<b>1261</b>	33	4	10	-	-	28	41	518
<b>VBD</b>	39	12	<b>1065</b>	-	-	14	-	-	-	-
<b>IN</b>	43	26	-	<b>870</b>	-	4	-	-	679	1
<b>RB</b>	411	51	2	<b>682</b>	<b>861</b>	137	-	-	-	1
<b>VBN</b>	<b>829</b>	26	-	-	-	22	-	<b>726</b>	-	-
<b>JJ</b>	-	664	495	21	51	<b>819</b>	2	32	372	125
<b>NNP</b>	333	-	17	1	11	321	<b>790</b>	4	28	15
<b>VBP</b>	22	9	15	1	7	194	-	31	4	-
<b>VBG</b>	322	30	-	-	-	462	-	-	1	-
<b>VBZ</b>	1	14	-	-	-	1	7	-	-	-



Aciertos: 608.633 (96,37 %)

Errores: 22.916 (3,63 %)

Se puede apreciar una leve mejoría en el porcentaje de etiquetas acertadas para el modelo que incorpora NFI; 96,25 % contra 96,47 % y 96,21 % contra 96,37 % para cada mitad respectivamente. Los errores más comunes son producidos en etiquetas NN del gold standard cuando son etiquetadas como JJ y NNP por TnT, para las dos mitades entrenadas tanto con WSJ como con WSJ + NFI. Se puede notar que el porcentaje de error al etiquetar JJ cuando era NN es menor en la evaluación realizada sobre TnT entrenado con el modelo que incorpora NFI.

A continuación se presentan las matrices de confusión entre las mitades de WSJ etiquetado con TnT entrenado con la mitad restante con y sin NFI.

**Cuadro 11:** 1 mitad WSJ etiquetado por TnT (entrenado con 2 mitad WSJ) vs 1 mitad WSJ etiquetado con TnT (entrenado con 2 mitad de WSJ + NFI)

TnT(2WSJ+NFI) TnT(2WSJ)	NN	JJ	VBN	NNP	VBD	VBG	VBP	VB	RP	IN
<b>JJ</b>	<b>592</b>	-	149	<b>353</b>	38	64	11	62	1	26
<b>NN</b>	-	<b>580</b>	30	<b>488</b>	16	<b>338</b>	97	<b>315</b>	-	3
<b>VBD</b>	8	37	<b>506</b>	13	-	-	6	7	-	-
<b>VBN</b>	17	<b>377</b>	-	7	<b>483</b>	-	4	8	-	-
<b>VB</b>	203	55	25	12	22	-	<b>321</b>	-	1	-
<b>RB</b>	29	147	1	14	1	-	-	2	292	229
<b>VBP</b>	99	3	-	1	11	-	-	223	-	2
<b>VBG</b>	167	217	-	36	1	-	-	1	-	-
<b>VBZ</b>	-	-	-	6	-	-	-	1	-	-
<b>NNS</b>	56	9	-	115	-	-	-	-	-	-

Aciertos: 621.391 (98,39 %)

Errores: 10.184 (1,61 %)

**Cuadro 12:** 2 mitad WSJ etiquetado por TnT (entrenado con 1 mitad WSJ) vs 2 mitad WSJ etiquetado con TnT (entrenado con 1 mitad de WSJ + NFI)

TnT(1WSJ+NFI) TnT(1WSJ)	JJ	VBN	NN	VBD	NNP	RP	IN	VB	VBP	VBG
<b>NN</b>	<b>765</b>	28	-	27	<b>436</b>	1	4	<b>292</b>	72	270
<b>VBD</b>	59	<b>525</b>	11	-	2	-	-	8	3	-
<b>JJ</b>	-	193	<b>497</b>	47	<b>335</b>	-	11	48	23	82
<b>VBN</b>	<b>286</b>	-	22	<b>439</b>	4	-	-	11	5	-
<b>RB</b>	160	-	38	-	30	<b>326</b>	<b>301</b>	15	3	-
<b>VB</b>	44	18	173	16	17	1	1	-	275	1
<b>VBZ</b>	-	-	-	-	5	-	-	1	-	-
<b>VBP</b>	14	3	101	7	4	-	1	219	-	-
<b>VBG</b>	190	1	156	2	18	-	-	6	2	-
<b>NNS</b>	9	-	50	-	72	-	-	-	-	-

Aciertos: 621.749 (98,45 %)

Errores: 9.802 (1,55 %)

La tercer evaluación de este experimento consiste en entrenar TnT con un cuarto de WSJ y con un cuarto de WSJ + NFI. Posteriormente con estos dos modelos se etiqueta los 3/4 restantes de WSJ y se construye la matriz de confusión. Se realiza la misma operación para cada uno de los cuartos:

**Cuadro 13:** *Rendimiento de TnT entrenado con cuartos de WSJ con y sin NFI*

Evaluación	Porcentaje de aciertos
TnT entrenado con el primer 1/4 de WSJ	95.93 %
TnT entrenado con el primer 1/4 de WSJ + NFI	96.26 %
TnT entrenado con el segundo 1/4 de WSJ	95.89 %
TnT entrenado con el segundo 1/4 de WSJ + NFI	96.26 %
TnT entrenado con el tercer 1/4 de WSJ	95.91 %
TnT entrenado con el tercer 1/4 de WSJ + NFI	96.29 %
TnT entrenado con el cuarto 1/4 de WSJ	95.9 %
TnT entrenado con el cuarto 1/4 de WSJ + NFI	96.30 %

En todos los casos se puede apreciar una mejora en el acierto de etiquetas para el modelo que incorpora NFI.

La cuarta evaluación de este experimento consiste en entrenar TnT con un décimo de WSJ y con un décimo de WSJ + NFI. Posteriormente con estos dos modelos se etiqueta los 9/10 restantes de WSJ y se presentan los resultados:

- 95.32 % de acierto de etiquetas para el etiquetado de 9/10 de WSJ con TnT entrenado con 1/10 WSJ
- 96.1 % de acierto de etiquetas para el etiquetado de 9/10 de WSJ con TnT entrenado con 1/10 WSJ+NFI

Se puede apreciar un aumento del porcentaje de aciertos en el modelo que incorpora NFI.

#### 4.2.2. Etiquetar el corpus WSJ con Stanford Tagger

La segunda evaluación de este experimento consiste en entrenar el etiquetador gramatical Stanford Tagger con WSJ como corpus de entrenamiento y con WSJ + NFI. Luego se procede a etiquetar el WSJ plano (sin etiquetas gramaticales) con estos dos modelos. Por último se contruye la matriz de confusión:

**Cuadro 14:** *WSJ original contra WSJ etiquetado con MaxEnt (entrenado con WSJ)*

MaxEnt(WSJ)	JJ	IN	NN	NNP	VBD	RB	VRN	VBP	RP	JJR
WSJ										
NN	<b>1726</b>	15	-	<b>1132</b>	16	61	18	28	1	2
RB	736	<b>1593</b>	189	139	-	-	3	1	293	36
JJ	-	60	<b>1276</b>	632	51	515	762	7	-	5
VRN	<b>894</b>	-	44	25	<b>1052</b>	1	-	5	-	-
NNPS	40	-	-	<b>997</b>	-	-	-	-	-	-
IN	87	-	4	22	-	<b>959</b>	-	2	527	-
VBG	196	-	<b>829</b>	14	-	-	1	1	-	-
VBD	40	-	26	8	-	-	<b>806</b>	14	-	-
RP	4	628	-	1	-	230	-	-	-	-
VB	58	8	365	36	37	12	22	544	-	6

Aciertos: 1.236.647 (97,90 %)

Errores: 26.477 (2,10 %)

**Cuadro 15:** *WSJ original contra WSJ etiquetado con MaxEnt (entrenado con WSJ + NFI)*

MaxEnt(WSJ+NFI)	JJ	IN	NNP	NN	RB	VBD	RP	VRN	VBP	NNPS
WSJ										
NN	<b>1918</b>	16	<b>1314</b>	-	73	21	3	19	30	-
RB	742	<b>1403</b>	141	177	-	-	527	3	3	-
JJ	-	68	685	<b>1260</b>	583	45	2	851	8	-
IN	107	-	29	3	<b>1062</b>	-	<b>1005</b>	-	2	-
VRN	<b>980</b>	-	29	49	-	<b>1049</b>	-	-	6	-
NNPS	39	-	<b>935</b>	-	-	-	-	-	-	-
VBD	34	-	10	24	-	-	-	<b>923</b>	15	-
VBG	294	-	25	817	-	-	-	1	1	-
NNP	555	29	-	458	21	4	1	9	4	546
VB	62	10	29	361	13	50	-	38	555	-

Aciertos: 1.234.495 (97,73 %)

Errores: 28.629 (2,27 %)

Se puede observar que el rendimiento del etiquetador entrenado con WSJ es un poco mejor (97,9 %) que cuando es entrenado con WSJ + NFI (97,73 %). La mayoría de los errores para Stanford Tagger entrenado con WSJ se da en

etiquetas NN del gold standard cuando son etiquetadas como JJ y NNP. Para Stanford Tagger entrenado con WSJ + NFI la mayoría de los errores se da en las mismas etiquetas, pero con cantidad de errores mayor, sobre todo para NN etiquetado como JJ.

La segunda evaluación de este experimento consiste en entrenar Stanford Tagger con la mitad de WSJ y con la mitad de WSJ + NFI. Posteriormente con estos dos modelos se etiqueta la mitad restante de WSJ y se construye la matriz de confusión. Se realiza la misma operación para cada mitad:

**Cuadro 16:** 1 mitad WSJ original contra 1 mitad WSJ etiquetado con MaxEnt (entrenado con 2 mitad de WSJ)

MaxEnt(2WSJ) WSJ	JJ	NN	NNP	IN	VBN	VBD	RB	NNS	VBG	JJR
NN	<b>1558</b>	-	<b>1027</b>	6	22	15	38	133	403	8
JJ	-	<b>1309</b>	<b>606</b>	32	<b>746</b>	39	299	65	263	5
RB	512	104	31	<b>989</b>	2	1	-	4	1	14
NNPS	31	-	<b>943</b>	-	-	-	-	246	-	-
VBN	545	28	36	-	-	<b>722</b>	-	-	-	-
VBG	192	<b>614</b>	22	-	1	-	-	-	-	-
VBD	41	26	9	-	<b>604</b>	-	-	-	-	-
NNP	401	542	-	23	8	10	38	156	37	-
IN	72	5	26	-	1	-	489	-	2	-
RP	2	3	1	449	-	-	179	-	-	-

Aciertos: 610.045 (96,59 %)

Errores: 21.529 (3,41 %)

**Cuadro 17:** 1 mitad WSJ original contra 1 mitad WSJ etiquetado con MaxEnt (entrenado con 2 mitad de WSJ + NFI)

MaxEnt(2WSJ+NFI) WSJ	JJ	NN	NNP	IN	VBN	VBD	RP	RB	NNS	VB
NN	<b>1434</b>	-	<b>1033</b>	6	13	14	1	41	124	124
JJ	-	<b>1119</b>	<b>571</b>	31	<b>625</b>	31	1	334	68	26
RB	470	81	73	<b>851</b>	2	1	251	-	1	19
NNPS	27	3	<b>834</b>	-	-	-	-	-	215	-
VBD	35	16	14	-	<b>748</b>	-	-	-	-	10
VBN	564	26	34	-	-	<b>642</b>	-	-	-	11
IN	80	7	27	-	-	-	<b>573</b>	531	-	1
VBG	262	531	26	-	1	-	-	-	-	-
NNP	445	497	-	19	3	6	-	24	141	16
VBZ	-	1	17	-	-	-	-	-	412	-

Aciertos: 611.099 (96,76 %)

Errores: 20.475 (3,24 %)

**Cuadro 18:** 2 mitad WSJ original contra 2 mitad WSJ etiquetado con MaxEnt (entrenado con 1 mitad de WSJ)

MaxEnt(1WSJ) WSJ	JJ	NN	IN	NNP	VBD	RB	VBN	NNPS	NNS	VBG
<b>NN</b>	<b>1604</b>	-	12	<b>916</b>	16	36	21	-	150	381
<b>JJ</b>	-	<b>1197</b>	45	522	37	381	483	-	46	202
<b>RB</b>	466	168	<b>944</b>	62	1	-	1	-	3	1
<b>VBN</b>	<b>863</b>	29	-	32	<b>779</b>	1	-	-	-	-
<b>IN</b>	50	3	-	26	-	<b>698</b>	-	-	-	2
<b>NNPS</b>	16	-	-	<b>651</b>	-	-	-	-	167	-
<b>VBG</b>	198	<b>572</b>	-	19	-	-	-	-	-	-
<b>VBD</b>	66	43	2	16	-	-	<b>570</b>	-	-	-
<b>NNP</b>	462	503	18	-	2	19	19	518	131	16
<b>RP</b>	3	1	426	-	-	129	-	-	-	-

Aciertos: 610.309 (96,64 %)

Errores: 21.241 (3,36 %)

**Cuadro 19:** 2 mitad WSJ original contra 2 mitad WSJ etiquetado con MaxEnt (entrenado con 1 mitad de WSJ)

MaxEnt(1WSJ+NFI) WSJ	JJ	NNP	NN	IN	VBN	RB	VBD	RP	NNPS	NNS
<b>NN</b>	<b>1482</b>	<b>1011</b>	-	10	18	42	12	2	-	146
<b>JJ</b>	-	522	<b>997</b>	43	444	382	26	8	1	40
<b>RB</b>	438	105	141	<b>819</b>	1	-	-	344	-	1
<b>VBN</b>	<b>810</b>	27	28	-	-	-	<b>706</b>	-	-	-
<b>VBD</b>	40	11	20	-	<b>742</b>	-	-	-	-	-
<b>IN</b>	50	29	3	-	-	<b>727</b>	-	<b>586</b>	-	-
<b>NNPS</b>	13	<b>597</b>	-	-	-	-	-	-	-	171
<b>VBG</b>	256	21	530	-	-	-	-	-	-	-
<b>NNP</b>	475	-	467	16	8	17	3	1	483	140
<b>VBZ</b>	-	9	1	-	-	-	-	-	2	406

Aciertos: 610.874 (96,73 %)

Errores: 20.676 (3,27 %)

Se puede apreciar una leve mejoría en el porcentaje de etiquetas acertadas; 96,23 % contra 96,46 % y 96,20 % contra 96,36 % para cada mitad respectivamente. Los errores más comunes son producidos en etiquetas NN del gold standard cuando son etiquetadas como JJ y NNP por TnT, para las dos mitades entrenadas tanto con WSJ como con WSJ + NFI. Se puede notar que el porcentaje de error al etiquetar JJ cuando era NN es menor en la evaluación realizada sobre TnT entrenado con WSJ + NFI.

A continuación se presentan las matrices de confusión entre las mitades de WSJ etiquetado con Stanford Tagger entrenado con la mitad restante con y sin NFI.

**Cuadro 20:** 1 mitad WSJ etiquetado por MaxEnt (entrenado con 2 mitad WSJ) vs 1 mitad WSJ etiquetado con MaxEnt (entrenado con 2 mitad de WSJ + NFI)

MaxEnt(2WSJ+NFI) MaxEnt(2WSJ)	JJ	NN	RP	VBN	RB	NNP	VB	IN	VBD	VBG
<b>NN</b>	<b>686</b>	-	2	13	46	<b>291</b>	<b>276</b>	-	21	213
<b>JJ</b>	-	<b>596</b>	1	175	202	183	48	15	19	49
<b>IN</b>	17	6	<b>507</b>	-	<b>318</b>	8	-	-	-	-
<b>VBD</b>	29	9	-	<b>460</b>	-	11	11	-	-	-
<b>VBN</b>	<b>309</b>	18	-	-	-	9	2	-	248	-
<b>NNP</b>	<b>268</b>	242	-	5	8	-	14	4	-	10
<b>WDT</b>	-	-	-	-	-	-	-	<b>252</b>	-	-
<b>VBP</b>	15	123	-	5	2	6	246	-	13	-
<b>RB</b>	128	17	206	1	-	65	14	141	-	-
<b>VBG</b>	199	196	-	-	1	26	-	-	-	-

Aciertos: 622.105 (98,50 %)

Errores: 9.469 (1,50 %)

**Cuadro 21:** 2 mitad WSJ etiquetado por MaxEnt (entrenado con 1 mitad WSJ) vs 2 mitad WSJ etiquetado con MaxEnt (entrenado con 1 mitad de WSJ + NFI)

MaxEnt(1WSJ+NFI) MaxEnt(1WSJ)	JJ	NN	VBN	RP	NNP	RB	VB	VBD	IN	NNS
<b>NN</b>	<b>611</b>	-	25	1	<b>344</b>	32	244	27	5	45
<b>JJ</b>	-	<b>513</b>	<b>254</b>	-	196	195	48	37	9	23
<b>VBD</b>	28	12	<b>494</b>	-	3	1	2	-	-	-
<b>IN</b>	9	3	-	<b>494</b>	18	<b>326</b>	-	2	-	-
<b>VBP</b>	17	132	6	-	4	-	<b>283</b>	15	2	1
<b>VBN</b>	246	18	-	-	7	-	9	<b>257</b>	-	-
<b>WDT</b>	-	-	-	-	-	-	-	-	<b>255</b>	-
<b>RB</b>	149	13	1	222	70	-	19	-	166	-
<b>NNP</b>	211	215	8	-	-	26	24	4	5	87
<b>VBZ</b>	-	1	-	-	5	1	-	-	-	208

Aciertos: 622.115 (98,51 %)

Errores: 9.435 (1,49 %)

La tercer evaluación de este experimento consiste en entrenar Stanford Tagger con un cuarto de WSJ y con un cuarto de WSJ + NFI. Posteriormente con estos dos modelos se etiqueta los 3/4 restantes de WSJ y se construye la matriz de confusión. Se realiza la misma operación para cada uno de los cuartos:

**Cuadro 22:** Rendimiento de Stanford Tagger entrenado con cuartos de WSJ con y sin NFI

Evaluación	Porcentaje de aciertos
------------	------------------------

**Cuadro 22:** Rendimiento de Stanford Tagger entrenado con cuartos de WSJ con y sin NFI

Evaluación	Porcentaje de aciertos
Stanford Tagger entrenado con el primer 1/4 de WSJ	96.30 %
Stanford Tagger entrenado con el primer 1/4 de WSJ + NFI	96.57 %
Stanford Tagger entrenado con el segundo 1/4 de WSJ	96.30 %
Stanford Tagger entrenado con el segundo 1/4 de WSJ + NFI	96.52 %
Stanford Tagger entrenado con el tercer 1/4 de WSJ	96.28 %
Stanford Tagger entrenado con el tercer 1/4 de WSJ + NFI	96.57 %
Stanford Tagger entrenado con el cuarto 1/4 de WSJ	96.24 %
Stanford Tagger entrenado con el cuarto 1/4 de WSJ + NFI	96.53 %

En todos los casos se puede apreciar una mejora en el acierto de etiquetas para el corpus de entrenamiento WSJ + NFI contra WSJ.

La cuarta evaluación de este experimento consiste en entrenar Stanford Tagger con un décimo de WSJ y con un décimo de WSJ + NFI. Posteriormente con estos dos modelos se etiqueta los 9/10 restantes de WSJ y se presentan los resultados:

- 95.67 % de acierto de etiquetas para el etiquetado de 9/10 de WSJ con Stanford Tagger entrenado con 1/10 WSJ
- 96.27 % de acierto de etiquetas para el etiquetado de 9/10 de WSJ con Stanford Tagger entrenado con 1/10 WSJ+NFI

Se puede apreciar un aumento del porcentaje de aciertos en el corpus de entrenamiento que incorpora NFI.

### 4.3. Tercer experimento: Etiquetar el corpus BNC

El tercer experimento realizado tiene como objetivo evaluar la nueva fuente de información obtenida (NFI) como corpus de entrenamiento. Para esto se entrenarán 2 etiquetadores gramaticales (Stanford Tagger y TnT) y se etiquetará con ellos el British National Corpus (BNC).

#### 4.3.1. Etiquetar el corpus BNC con TnT

La primer evaluación de este experimento consiste en entrenar el etiquetador gramatical TnT con WSJ como corpus de entrenamiento y con WSJ + NFI. Luego se procede a etiquetar el BNC plano (sin etiquetas gramaticales) con estos dos modelos. Por último se contruye la matriz de confusión:

**Cuadro 23:** *BNC original contra BNC etiquetado con TnT (entrenado con WSJ)*

TnT(BNC) BNC	NNP	JJ	NN	VBN	NNS	WRB	CD	NNPS	VBD	RB
<b>NN1</b>	<b>26585</b>	<b>5739</b>	-	146	732	3	52	13	157	378
<b>AJ0</b>	<b>8608</b>	-	2352	<b>3680</b>	50	1	34	24	327	1092
<b>DT0</b>	83	<b>7771</b>	208	-	1	-	-	-	-	988
<b>AV0</b>	1014	2028	<b>4277</b>	10	242	6	188	-	7	-
<b>NN0</b>	469	342	-	7	<b>3133</b>	-	<b>2601</b>	18	2	8
<b>CJS</b>	217	315	581	20	76	<b>2903</b>	-	-	43	1202
<b>NN2</b>	2472	75	690	2	-	-	77	<b>2578</b>	-	5
<b>VVN</b>	56	361	94	-	-	-	-	-	2365	1
<b>VVD</b>	54	235	87	2159	-	-	-	-	-	6
<b>AVP</b>	25	7	140	-	-	1	-	-	-	2070

Aciertos: 1.849.040 (92,47 %)

Errores: 150.675 (7,53 %)

**Cuadro 24:** *BNC original contra BNC etiquetado con TnT (entrenado con WSJ + NFI)*

TnT(WSJ+NFI) BNC	NNP	JJ	NN	NNS	VBN	WRB	NNPS	VBD	CD	WDT
<b>NN1</b>	<b>26413</b>	<b>4240</b>	-	663	101	1	8	118	56	1
<b>AJ0</b>	<b>8621</b>	-	1806	33	<b>3361</b>	1	21	291	20	1
<b>DT0</b>	75	<b>7689</b>	206	1	-	-	-	-	-	802
<b>AV0</b>	1068	1756	<b>4347</b>	226	16	2	-	4	125	4
<b>NN0</b>	447	474	-	<b>3366</b>	9	-	14	-	2120	-
<b>CJS</b>	222	252	705	101	18	<b>2903</b>	6	45	-	54
<b>NN2</b>	2328	85	795	-	-	-	<b>2542</b>	-	19	-
<b>VVN</b>	46	458	79	-	-	-	-	<b>2349</b>	-	-
<b>VVD</b>	53	217	76	-	2343	-	-	-	-	-
<b>CJT</b>	-	-	1	-	-	-	-	-	-	1857



Aciertos: 1.854.577 (92,74 %)

Errores: 145.138 (7,26 %)

Se puede observar que el rendimiento del etiquetador TnT entrenado con WSJ+NFI es un poco mejor (97,14 %) que el rendimiento de TnT entrenado con WSJ (97,1 %). La mayoría de los errores para TnT entrenado con WSJ se da en etiquetas NN del gold standard cuando son etiquetadas como JJ y NNP por TnT. Para TnT entrenado con WSJ + NFI la mayoría de los errores se da en las mismas etiquetas, pero con cantidad de errores mayor, sobre todo para NN etiquetado como NNP.

La segunda evaluación de este experimento consiste en entrenar TnT con la mitad de WSJ y con la mitad de WSJ + NFI. Posteriormente con estos dos modelos se etiqueta BNC y se construye la matriz de confusión. Se realiza la misma operación para cada mitad:

**Cuadro 25:** BNC original contra BNC etiquetado con TnT (entrenado con 2 mitad de WSJ)

TnT(2WSJ) BNC	NNP	JJ	VBN	NN	NNS	WRB	VBD	RB	VBG	NNPS
NN1	<b>27194</b>	<b>6142</b>	180	-	804	5	191	396	1735	9
AJ0	<b>8813</b>	-	<b>4421</b>	<b>3008</b>	69	1	371	1106	2191	22
DT0	82	<b>7802</b>	-	200	1	-	-	903	-	-
AV0	1008	2145	15	<b>4196</b>	314	5	6	-	41	2
NN0	483	650	8	-	<b>3410</b>	-	1	9	4	19
NN2	<b>2962</b>	76	1	734	-	-	-	19	-	2174
CJS	202	314	27	566	86	<b>2904</b>	45	1152	15	-
VVN	44	432	-	106	-	-	2707	1	7	-
VVD	45	258	2390	122	-	-	-	7	20	-
AVP	18	7	-	144	-	-	-	2335	-	-

Aciertos: 1.841.617 (92,09 %)

Errores: 158.098 (7,91 %)

**Cuadro 26:** BNC original contra BNC etiquetado con TnT (entrenado con 2 mitad de WSJ+NFI)

TnT(2WSJ+NFI) BNC	NNP	JJ	NN	VBN	NNS	WRB	NNPS	VBD	CD	WDT
NN1	<b>26728</b>	<b>4264</b>	-	101	663	1	8	119	61	1
AJ0	<b>8795</b>	-	1793	<b>3481</b>	34	-	21	295	22	-
DT0	86	<b>7648</b>	204	-	1	-	-	-	-	869
AV0	1074	1881	<b>4382</b>	15	216	1	-	4	130	4
NN0	445	424	-	9	<b>3380</b>	-	20	-	2177	-
CJS	229	260	738	15	102	<b>2903</b>	1	48	-	17
NN2	<b>2670</b>	76	795	-	-	-	2331	-	22	-
VVD	48	227	87	<b>2458</b>	-	-	-	-	-	-

**Cuadro 26:** BNC original contra BNC etiquetado con TnT (entrenado con 2 mitad de WSJ+NFI)

$\begin{matrix} \text{TnT(2WSJ+NFI)} \\ \text{BNC} \end{matrix}$	NNP	JJ	NN	VBN	NNS	WRB	NNPS	VBD	CD	WDT
<b>VVN</b>	41	481	82	-	-	-	-	2318	-	-
<b>CJT</b>	6	-	1	-	-	-	-	-	-	1857

Aciertos: 1.853.072 (92,67 %)

Errores: 146.643 (7,33 %)

**Cuadro 27:** BNC original contra BNC etiquetado con TnT (entrenado con 1 mitad de WSJ)

$\begin{matrix} \text{TnT(1WSJ)} \\ \text{BNC} \end{matrix}$	NNP	JJ	NN	VBN	NNS	WRB	NNPS	VBD	VBG	RB
<b>NN1</b>	<b>26286</b>	<b>6528</b>	-	187	783	3	36	169	1614	435
<b>AJ0</b>	<b>8503</b>	-	<b>2921</b>	<b>3713</b>	64	1	28	476	2027	1238
<b>DT0</b>	84	<b>7763</b>	232	-	1	-	-	-	-	976
<b>AV0</b>	1157	2031	<b>4455</b>	30	515	1	1	8	42	-
<b>NN0</b>	478	796	-	6	<b>3225</b>	-	21	2	14	4
<b>CJS</b>	193	235	626	24	77	<b>2903</b>	6	51	14	1165
<b>NN2</b>	2504	87	863	2	-	-	<b>2714</b>	-	-	1
<b>VVN</b>	57	615	95	-	-	-	-	2581	13	5
<b>PRP</b>	733	1415	2260	57	499	3	2	-	467	614
<b>VVD</b>	67	317	99	2126	1	-	-	-	27	5

Aciertos: 1.842.527 (92,14 %)

Errores: 157.188 (7,86 %)

**Cuadro 28:** BNC original contra BNC etiquetado con TnT (entrenado con 1 mitad de WSJ+NFI)

$\begin{matrix} \text{TnT(1WSJ+NFI)} \\ \text{BNC} \end{matrix}$	NNP	JJ	NN	NNS	VBN	WRB	NNPS	VBD	CD	WDT
<b>NN1</b>	<b>26360</b>	<b>4345</b>	-	677	107	2	14	127	63	3
<b>AJ0</b>	<b>8729</b>	-	1808	35	<b>3290</b>	1	20	303	24	1
<b>DT0</b>	91	<b>7655</b>	210	1	-	-	-	-	-	870
<b>AV0</b>	1225	1630	<b>4409</b>	215	18	-	1	4	143	4
<b>NN0</b>	442	501	-	<b>3321</b>	8	-	21	-	2174	-
<b>CJS</b>	221	256	696	100	18	<b>2903</b>	6	47	-	63
<b>NN2</b>	2327	85	815	-	-	-	<b>2568</b>	-	17	-
<b>VVD</b>	60	218	76	-	<b>2429</b>	-	-	-	-	-
<b>VVN</b>	62	481	84	-	-	-	-	2381	-	-
<b>CJT</b>	-	-	1	-	-	-	-	-	-	1867

Aciertos: 1.853.701 (92,70 %)

Errores: 146.014 (7,30 %)

Se puede apreciar una leve mejoría en el porcentaje de etiquetas acertadas para el modelo que incorpora NFI; 92,09 % contra 92,67 % y 92,14 % contra 92,7 % para cada mitad respectivamente. Los errores más comunes son producidos en etiquetas NN del gold standard cuando son etiquetadas como JJ y NNP por TnT, para las dos mitades entrenadas tanto con WSJ como con WSJ + NFI. Se puede notar que el porcentaje de error al etiquetar JJ cuando era NN es menor en la evaluación realizada sobre TnT entrenado con el modelo que incorpora NFI.

A continuación se presentan las matrices de confusión entre las mitades de WSJ etiquetado con TnT entrenado con la mitad restante con y sin NFI.

**Cuadro 29:** BNC etiquetado por TnT (entrenado con 1 mitad WSJ) vs BNC etiquetado con TnT (entrenado con 1 mitad de WSJ + NFI)

TnT(1WSJ+NFI) TnT(1WSJ)	NN	JJ	VBN	VB	NNP	NNS	VBD	VBP	VBG	RB
<b>JJ</b>	<b>5144</b>	-	1092	693	<b>1701</b>	128	283	82	474	913
<b>NN</b>	-	<b>3630</b>	65	<b>2399</b>	<b>2215</b>	390	58	417	1168	661
<b>VBD</b>	115	366	<b>2487</b>	93	47	1	-	21	-	7
<b>VB</b>	<b>1960</b>	291	131	-	230	1	97	1290	5	103
<b>VBZ</b>	12	65	-	5	38	<b>1931</b>	4	3	-	16
<b>NNP</b>	<b>1870</b>	1178	42	270	-	618	21	32	87	152
<b>VBN</b>	155	1512	-	70	78	1	<b>1865</b>	16	5	12
<b>VBP</b>	641	125	24	1343	41	2	38	-	-	27
<b>RB</b>	491	1215	4	238	279	5	3	9	-	-
<b>VBG</b>	1028	1198	5	21	188	5	23	7	-	5

Aciertos: 1.938.000 (96,91 %)

Errores: 61.726 (3,09 %)

**Cuadro 30:** BNC etiquetado por TnT (entrenado con 2 mitad WSJ) vs BNC etiquetado con TnT (entrenado con 2 mitad de WSJ + NFI)

TnT(2WSJ+NFI) TnT(2WSJ)	NN	JJ	VBN	VB	NNP	NNS	VBD	VBP	RP	VBG
<b>JJ</b>	<b>5020</b>	-	846	433	1470	73	260	77	2	324
<b>NN</b>	-	<b>4138</b>	67	<b>2355</b>	<b>2027</b>	180	81	455	3	932
<b>VBD</b>	133	304	<b>2442</b>	46	38	2	-	25	-	-
<b>NNP</b>	<b>2220</b>	1337	23	242	-	708	10	22	-	85
<b>VB</b>	<b>2193</b>	359	116	-	140	3	189	1602	2	4
<b>VBN</b>	185	<b>2091</b>	-	53	28	-	<b>1882</b>	18	-	2
<b>VBZ</b>	16	31	-	3	51	<b>1998</b>	-	10	-	-
<b>VBP</b>	599	50	10	1271	37	1	69	-	-	-
<b>RB</b>	539	1146	-	467	237	26	1	24	1177	-
<b>VBG</b>	1166	1172	-	5	200	3	10	-	-	-

Aciertos: 1.938.152 (96,92 %)

Errores: 61.574 (3,08 %)

La tercer evaluación de este experimento consiste en entrenar TnT con un cuarto de WSJ y con un cuarto de WSJ + NFI. Posteriormente con estos dos modelos se etiqueta los 3/4 restantes de WSJ y se construye la matriz de confusión. Se realiza la misma operación para cada uno de los cuartos:

**Cuadro 31:** Rendimiento de TnT entrenado con cuartos de WSJ con y sin NFI

Evaluación	Porcentaje de aciertos
TnT entrenado con el primer 1/4 de WSJ	91.75 %
TnT entrenado con el primer 1/4 de WSJ + NFI	92.62 %
TnT entrenado con el segundo 1/4 de WSJ	91.74 %
TnT entrenado con el segundo 1/4 de WSJ + NFI	92.63 %
TnT entrenado con el tercer 1/4 de WSJ	91.64 %
TnT entrenado con el tercer 1/4 de WSJ + NFI	92.62 %
TnT entrenado con el cuarto 1/4 de WSJ	91.64 %
TnT entrenado con el cuarto 1/4 de WSJ + NFI	92.58 %

En todos los casos se puede apreciar una mejora en el acierto de etiquetas para el modelo que incorpora NFI.

La cuarta evaluación de este experimento consiste en entrenar TnT con un décimo de WSJ y con un décimo de WSJ + NFI. Posteriormente con estos dos modelos se etiqueta BNC y se presentan los resultados:

- 90.9 % de acierto de etiquetas para el etiquetado de BNC con TnT entrenado con 1/10 WSJ
- 92.55 % de acierto de etiquetas para el etiquetado de BNC con TnT entrenado con 1/10 WSJ+NFI

Se puede apreciar un aumento del porcentaje de aciertos en el modelo que incorpora NFI.

#### 4.3.2. Etiquetar el corpus BNC con Stanford Tagger

La segunda evaluación de este experimento consiste en entrenar el etiquetador gramatical Stanford Tagger con WSJ como corpus de entrenamiento y con WSJ + NFI. Luego se procede a etiquetar el BNC plano (sin etiquetas gramaticales) con estos dos modelos. Por último se contruye la matriz de confusión:

**Cuadro 32:** *BNC original contra BNC etiquetado con MaxEnt (entrenado con WSJ)*

MaxEnt(BNC) BNC	NNP	JJ	NN	VBN	NNS	WRB	RB	CD	VBG	NNPS
NN1	<b>26141</b>	<b>4045</b>	-	115	533	-	253	16	1143	7
AJ0	<b>8675</b>	-	<b>2860</b>	<b>3276</b>	30	-	1033	3	2054	12
DT0	119	<b>8132</b>	192	-	5	-	443	-	-	-
AV0	982	2314	<b>3753</b>	159	236	-	-	160	85	4
NN0	567	1115	-	19	<b>3132</b>	-	10	2605	2	5
NN2	<b>2982</b>	90	750	-	-	-	6	47	-	1959
CJS	60	334	247	57	104	<b>2901</b>	522	1	35	2
AVP	43	17	137	-	-	-	2691	-	-	-
UNC	1754	255	540	9	199	-	2	426	1	18
CJT	-	1	-	-	-	-	-	-	-	-

Aciertos: 1.856.979 (92,86 %)

Errores: 142.739 (7,14 %)

**Cuadro 33:** *BNC original contra BNC etiquetado con MaxEnt (entrenado con WSJ + NFI)*

MaxEnt(WSJ+NFI) BNC	NNP	JJ	NN	CD	NNS	VBN	WRB	RB	NNPS	VBG
NN1	<b>26206</b>	<b>3864</b>	-	22	663	108	-	277	1	1166
AJ0	<b>8263</b>	-	2099	4	25	<b>3145</b>	-	876	12	1707
DT0	109	<b>7836</b>	188	-	4	-	-	793	-	-
AV0	840	2210	<b>3640</b>	195	279	123	-	-	2	70
NN0	469	863	-	<b>3222</b>	<b>3146</b>	6	-	-	8	9
CJS	102	374	357	-	147	37	<b>2901</b>	776	2	39
NN2	<b>2643</b>	68	783	74	-	1	-	8	1844	-
AVP	39	6	140	-	-	-	-	2101	-	-
PRP	497	1680	887	1	572	115	-	711	1	624
VVN	76	461	87	-	1	-	-	1	-	7

Aciertos: 1.859.888 (93,01 %)

Errores: 139.830 (6,99 %)

Se puede observar que el rendimiento del etiquetador entrenado con WSJ es un poco mejor (93,01 %) que cuando es entrenado con WSJ + NFI (92,86 %). La mayoría de los errores para Stanford Tagger entrenado con WSJ se da en

etiquetas NN del gold standard cuando son etiquetadas como JJ y NNP. Para Stanford Tagger entrenado con WSJ + NFI la mayoría de los errores se da en las mismas etiquetas, pero con cantidad de errores mayor, sobre todo para NN etiquetado como JJ.

La segunda evaluación de este experimento consiste en entrenar Stanford Tagger con la mitad de WSJ y con la mitad de WSJ + NFI. Posteriormente con estos dos modelos se etiqueta BNC y se construye la matriz de confusión. Se realiza la misma operación para cada mitad:

**Cuadro 34:** *BNC original contra BNC etiquetado con MaxEnt (entrenado con 2 mitad de WSJ)*

MaxEnt(2WSJ) BNC	NNP	JJ	NN	NNS	VBN	WRB	RB	CD	VBG	NNPS
<b>NN1</b>	<b>26061</b>	<b>4689</b>	-	620	97	-	248	6	1296	8
<b>AJ0</b>	<b>8570</b>	-	<b>3574</b>	49	<b>3001</b>	-	1100	6	2190	10
<b>DT0</b>	133	<b>8068</b>	224	1	-	-	478	1	-	-
<b>AV0</b>	985	2428	<b>3611</b>	475	143	-	-	164	77	9
<b>NN0</b>	554	1404	-	<b>3128</b>	7	-	14	2633	7	6
<b>NN2</b>	<b>2984</b>	146	855	-	-	-	4	46	-	2054
<b>CJS</b>	98	210	237	146	34	<b>2901</b>	568	3	21	2
<b>AVP</b>	47	6	152	-	-	-	2793	-	-	-
<b>UNC</b>	1763	268	454	212	3	-	3	518	1	20
<b>CJT</b>	-	1	-	-	-	-	-	-	-	-

Aciertos: 1.851.792 (92,60 %)

Errores: 147.926 (7,40 %)

**Cuadro 35:** *BNC original contra BNC etiquetado con MaxEnt (entrenado con 2 mitad de WSJ+NFI)*

MaxEnt(2WSJ+NFI) BNC	NNP	JJ	NN	CD	NNS	VBN	WRB	RB	NNPS	VBG
<b>NN1</b>	<b>26036</b>	<b>3867</b>	-	17	657	109	-	285	3	1179
<b>AJ0</b>	<b>8130</b>	-	2151	5	24	<b>3011</b>	-	897	13	1653
<b>DT0</b>	108	<b>7742</b>	215	-	3	-	-	867	-	-
<b>AV0</b>	875	2220	<b>3538</b>	203	276	129	-	-	1	68
<b>NN0</b>	452	866	-	<b>3240</b>	<b>3138</b>	8	-	-	11	8
<b>CJS</b>	123	327	333	2	158	58	<b>2901</b>	908	2	17
<b>NN2</b>	<b>2513</b>	75	800	81	-	1	-	10	1922	-
<b>AVP</b>	40	5	142	-	-	-	-	2014	-	-
<b>VVD</b>	75	203	89	-	-	1573	-	9	-	13
<b>PRP</b>	495	1528	758	2	689	109	-	723	-	587

Aciertos: 1.859.947 (93,01 %)

Errores: 139.771 (6,99 %)

**Cuadro 36:** BNC original contra BNC etiquetado con MaxEnt (entrenado con 1 mitad de WSJ)

MaxEnt(1WSJ) BNC	NNP	JJ	NN	VBN	NNS	WRB	RB	CD	VBG	VBD
NN1	<b>27101</b>	<b>4519</b>	-	146	550	-	241	7	1369	116
AJ0	<b>9043</b>	-	<b>3838</b>	<b>3837</b>	43	-	1025	3	2238	296
DT0	128	<b>8324</b>	185	-	1	-	461	-	-	-
AV0	1071	2583	<b>3445</b>	141	311	-	-	170	71	78
NN2	<b>3415</b>	83	885	-	-	-	7	41	-	1
NN0	573	955	-	24	<b>3189</b>	-	5	2732	3	8
CJS	82	412	267	54	106	<b>2901</b>	400	1	23	63
AVP	21	26	140	-	-	-	2859	-	-	-
VVN	85	464	126	-	1	-	1	-	1	1986
UNC	1757	181	509	9	242	-	3	442	1	5

Aciertos: 1.848.799 (92,45 %)

Errores: 150.919 (7,55 %)

**Cuadro 37:** BNC original contra BNC etiquetado con MaxEnt (entrenado con 1 mitad de WSJ+NFI)

MaxEnt(1WSJ+NFI) BNC	NNP	JJ	NN	CD	VBN	NNS	WRB	RB	VBD	NNPS
NN1	<b>26776</b>	<b>3786</b>	-	20	109	684	-	282	94	1
AJ0	<b>8368</b>	-	2113	3	<b>3249</b>	32	-	853	215	10
DT0	108	<b>7768</b>	192	-	-	1	-	883	-	-
AV0	950	2360	<b>3475</b>	194	90	325	-	-	39	2
NN0	473	752	-	<b>3302</b>	2	<b>3193</b>	-	-	4	4
CJS	151	409	348	-	33	150	<b>2901</b>	856	39	1
NN2	<b>2831</b>	66	797	74	1	-	-	8	1	1636
AVP	32	6	139	-	-	-	-	2080	-	-
VVN	77	492	94	-	-	1	-	1	1756	-
PRP	605	1613	925	-	129	741	-	802	122	-

Aciertos: 1.857.971 (92,91 %)

Errores: 141.747 (7,09 %)

Se puede apreciar una leve mejoría en el porcentaje de etiquetas acertadas; 92,6 % contra 93,01 % y 92,45 % contra 92,91 % para cada modelo respectivamente. Los errores más comunes son producidos en etiquetas NN del gold standard cuando son etiquetadas como JJ y NNP por , para las dos mitades entrenadas tanto con WSJ como con WSJ + NFI. Se puede notar que el porcentaje de error al etiquetar JJ cuando era NN es menor en la evaluación realizada sobre Stanford Tagger entrenado con WSJ + NFI.

A continuación se presentan las matrices de confusión para BNC etiquetado con Stanford Tagger entrenado con la mitad de WSJ con y sin NFI.

**Cuadro 38:** BNC etiquetado por MaxEnt (entrenado con 2 mitad WSJ) vs BNC etiquetado con MaxEnt (entrenado con 2 mitad de WSJ + NFI)

MaxEnt(2WSJ+NFI) MaxEnt(2WSJ)	JJ	NN	NNP	VBN	NNS	VB	RP	VBG	RB	VBD
NN	<b>4691</b>	-	<b>2413</b>	152	392	<b>1859</b>	11	1579	813	192
JJ	-	<b>3864</b>	1168	1524	85	356	2	336	1233	201
NNP	<b>1688</b>	<b>2727</b>	-	80	938	308	-	117	267	61
VBD	310	107	17	<b>2345</b>	-	135	-	6	21	-
VBZ	23	16	35	1	<b>1922</b>	6	-	-	24	9
IN	146	247	203	17	32	443	<b>1854</b>	75	1408	102
VBP	190	1075	35	28	15	<b>1693</b>	-	4	64	123
VBN	1351	99	69	-	13	66	-	12	20	1362
VBG	1220	1166	139	11	-	25	-	-	3	2
VB	300	1125	295	135	11	-	-	17	127	88

Aciertos: 1.933.574 (96,69 %)

Errores: 66.144 (3,31 %)

**Cuadro 39:** BNC etiquetado por MaxEnt (entrenado con 1 mitad WSJ) vs BNC etiquetado con MaxEnt (entrenado con 1 mitad de WSJ + NFI)

MaxEnt(1WSJ+NFI) MaxEnt(1WSJ)	JJ	NN	NNP	VBN	NNS	VB	RP	RB	VBG	VBD
NN	<b>5058</b>	-	<b>2370</b>	127	414	<b>1895</b>	15	372	1438	115
JJ	-	<b>3812</b>	1138	1148	80	360	8	1331	272	130
NNP	<b>2055</b>	<b>2670</b>	-	44	957	280	-	218	104	24
VBD	384	188	31	<b>2361</b>	-	186	-	21	28	-
VBZ	20	63	35	-	<b>2074</b>	37	-	43	1	12
VBN	<b>1973</b>	166	73	-	5	62	-	30	15	1366
IN	169	302	246	14	49	344	<b>1807</b>	1485	63	18
VBP	232	1021	27	27	12	1510	-	77	5	150
RB	1218	453	209	13	34	241	1081	-	21	13
VB	302	1204	269	145	20	-	-	139	23	107

Aciertos: 1.933.672 (96,70 %)

Errores: 66.046 (3,30 %)

La tercer evaluación de este experimento consiste en entrenar Stanford Tagger con un cuarto de WSJ y con un cuarto de WSJ + NFI. Posteriormente con estos dos modelos se etiqueta BNC y se construye la matriz de confusión. Se realiza la misma operación para cada uno de los cuartos:

**Cuadro 40:** Rendimiento de Stanford Tagger entrenado con cuartos de WSJ con y sin NFI

Evaluación	Porcentaje de aciertos
------------	------------------------



**Cuadro 40:** Rendimiento de Stanford Tagger entrenado con cuartos de WSJ con y sin NFI

Evaluación	Porcentaje de aciertos
Stanford Tagger entrenado con el primer 1/4 de WSJ	92.09 %
Stanford Tagger entrenado con el primer 1/4 de WSJ + NFI	92.92 %
Stanford Tagger entrenado con el segundo 1/4 de WSJ	92.10 %
Stanford Tagger entrenado con el segundo 1/4 de WSJ + NFI	92.91 %
Stanford Tagger entrenado con el tercer 1/4 de WSJ	92.14 %
Stanford Tagger entrenado con el tercer 1/4 de WSJ + NFI	92.89 %
Stanford Tagger entrenado con el cuarto 1/4 de WSJ	91.98 %
Stanford Tagger entrenado con el cuarto 1/4 de WSJ + NFI	92.83 %

En todos los casos se puede apreciar una mejora en el acierto de etiquetas para el corpus de entrenamiento WSJ + NFI contra WSJ.

La cuarta evaluación de este experimento consiste en entrenar Stanford Tagger con un décimo de WSJ y con un décimo de WSJ + NFI. Posteriormente con estos dos modelos se etiqueta BNC y se presentan los resultados:

- 91.25 % de acierto de etiquetas para el etiquetado de BNC con Stanford Tagger entrenado con 1/10 WSJ
- 92.81 % de acierto de etiquetas para el etiquetado de BNC con Stanford Tagger entrenado con 1/10 WSJ+NFI

Se puede apreciar un aumento del porcentaje de aciertos en el corpus de entrenamiento que incorpora NFI.

## 5. Conclusiones

## Referencias

- [1] Jurafsky, D. Martin, J. H., Speech and Language Processing: An introduction to natural language processing, computational linguistics, and speech recognition, Second edition, chapter 5, New Jersey: Prentice Hall.
- [2] Chris Manning and Hinrich Schütze, Foundations of Statistical Natural Language Processing, MIT Press. Cambridge, MA: May 1999
- [3] Thorsten Brants, TnT: a statistical part-of-speech tagger, Proceedings of the sixth conference on Applied natural language processing, p.224-231, April 29-May 04, 2000, Seattle, Washington
- [4] Kristina Toutanova and Christopher D. Manning. 2000. Enriching the Knowledge Sources Used in a Maximum Entropy Part-of-Speech Tagger. In Proceedings of the Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora (EMNLP/VLC-2000), pp. 63-70.
- [5] Mitchell P. Marcus , Mary Ann Marcinkiewicz , Beatrice Santorini, Building a large annotated corpus of English: the penn treebank, Computational Linguistics, v.19 n.2, June 1993
- [6] Reference Guide for the British National Corpus (World Edition) edited by Lou Burnard, October 2000
- [7] Stevenson M., A corpus-based approach to deriving lexical mappings, EACL '99 Proceedings of the ninth conference on European chapter of the Association for Computational Linguistics, Pages 285-286
- [8] Brown K. (Editor) 2005. Encyclopedia of Language and Linguistics 2nd Edition. Oxford: Elsevier.
- [9] Sinclair, J. 'The automatic analysis of corpora', in Svartvik, J. (ed.) Directions in Corpus Linguistics (Proceedings of Nobel Symposium 82). Berlin: Mouton de Gruyter. 1992.