

UNIVERSIDAD DE BUENOS AIRES  
Facultad de Ciencias Exactas y Naturales  
Departamento de Computación

## **Nuevas fuentes de información para entrenamiento de etiquetadores gramaticales**

Tesista: Fernando Jorge Rodriguez  
Director: Dr. José Castaño

Buenos Aires, Marzo de 2012.



## 0.1. Experimentación

### 0.1.1. 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	are
cat NN	often
is	kept
a	as
small	pets
furry	.
animal	She
with	put
a	out
tail	a
,	hand
whiskers	and
,	stroked
and	the
sharp	cat NN
claws	softly
that	...
kills	...
smaller	domestic
animals	animals
such	such
as	as
mice	dogs
and	and
birds	cats NNS
.	.
Cats NNS	

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>1</sup> y luego se realizará la comparación. La matriz de confusión<sup>2</sup> generada a partir de dicha comparación es la siguiente:

**Cuadro 1:** 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	-

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.

#### 0.1.2. Segundo experimento: entrenamiento de TnT con la nueva fuente de información generada

El segundo experimento realizado tiene como objetivo evaluar la nueva fuente de información obtenida (NFI) como corpus de entrenamiento. Para esto se utilizará el Wall Street Journal (WSJ), parte de Penn Tree Bank, como corpus objetivo.

La primer evaluación de este segundo experimento consiste en entrenar el etiquetador gramatical con WSJ como corpus de entrenamiento y con WSJ +

<sup>1</sup>Wall Street Journal es un corpus anotado, parte del Penn Treebank

<sup>2</sup>Las matrices de confusión presentadas de aquí en adelante contienen las primeras 10 etiquetas de mayor error

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 2:** *WSJ original contra WSJ etiquetado con TnT (entrenado con WSJ)*

TnT(WSJ) WSJ	JJ	NN	NNP	VBN	VBD	IN	RB	RP	NNPS	VBG
<b>JJ</b>	-	909	717	<b>932</b>	51	71	569	23	1	229
<b>NN</b>	<b>2537</b>	-	<b>1680</b>	30	40	19	72	7	-	622
<b>NNP</b>	454	358	-	4	3	26	47	1	<b>1262</b>	7
<b>VBN</b>	<b>1215</b>	37	22	-	<b>1151</b>	-	-	-	-	-
<b>VBD</b>	60	40	6	<b>1522</b>	-	-	-	-	-	-
<b>IN</b>	109	4	28	-	-	-	<b>1429</b>	<b>1353</b>	-	2
<b>RB</b>	752	203	46	4	1	<b>1500</b>	-	842	-	1
<b>RP</b>	2	-	1	-	-	371	144	-	-	-
<b>NNPS</b>	34	-	418	-	-	-	-	-	-	-
<b>VBG</b>	411	879	19	-	-	-	-	-	-	-

Porcentaje de aciertos: 97,38 %

**Cuadro 3:** *WSJ original contra WSJ etiquetado con TnT (entrenado con WSJ + NFI)*

TnT(WSJ+NFI) WSJ	JJ	NN	NNP	VBN	VBD	IN	RB	RP	NNPS	VB
<b>JJ</b>	-	895	1032	908	53	83	620	28	1	77
<b>NN</b>	<b>2776</b>	-	<b>2087</b>	48	43	16	83	7	-	337
<b>NNP</b>	361	256	-	4	2	19	45	1	<b>1255</b>	11
<b>VBN</b>	<b>1326</b>	35	27	-	<b>1236</b>	-	-	-	-	43
<b>VBD</b>	69	24	10	<b>1760</b>	-	-	-	-	-	42
<b>IN</b>	107	4	46	-	-	-	<b>1163</b>	<b>1584</b>	-	3
<b>RB</b>	834	193	63	4	1	<b>1705</b>	-	<b>1138</b>	-	50
<b>RP</b>	2	1	1	-	-	348	88	-	-	-
<b>NNPS</b>	33	-	448	-	-	-	-	-	-	-
<b>VB</b>	89	408	34	48	42	13	11	-	-	-

Porcentaje de aciertos: 97,12 %

Se puede observar que el rendimiento del etiquetador TnT entrenado con WSJ es un poco mejor (97,38 %) que el rendimiento de TnT entrenado con WSJ + NFI (97,12 %). 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 JJ.

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 4:** *Primer mitad WSJ original vs primer mitad WSJ etiquetado con TnT (entrenado con la segunda mitad de WSJ)*

TnT(2WSJ) WSJ	JJ	NN	NNP	VBN	VBD	IN	RB	VB	VBP	RP
JJ	-	<b>1042</b>	547	<b>810</b>	54	20	312	61	6	6
NN	<b>1968</b>	-	<b>1167</b>	25	24	6	55	272	72	2
NNP	422	520	-	9	11	19	43	21	3	-
VBN	605	25	24	-	<b>831</b>	-	-	38	2	-
VBD	75	21	12	<b>1140</b>	-	-	1	33	8	-
IN	70	4	24	1	-	-	<b>617</b>	1	4	<b>616</b>
RB	437	91	26	3	2	<b>844</b>	-	19	3	316
VB	77	414	41	25	13	6	8	-	343	-
VBP	26	296	19	19	33	5	4	<b>623</b>	-	-
RP	-	-	2	-	-	230	122	-	-	-

Porcentaje de aciertos: 96,23 %

**Cuadro 5:** *Primer mitad WSJ original vs primer mitad WSJ etiquetado con TnT (entrenado con segunda mitad de WSJ + NFI)*

TnT(2WSJ+NFI) WSJ	JJ	NN	NNP	VBN	VBD	IN	RB	RP	VBG	NNPS
JJ	-	<b>879</b>	<b>674</b>	<b>616</b>	29	37	337	7	210	-
NN	<b>1742</b>	-	<b>1293</b>	29	27	6	65	2	560	1
NNP	336	354	-	6	4	15	42	-	22	557
VBN	<b>636</b>	15	25	-	<b>716</b>	-	-	-	-	-
VBD	56	10	17	<b>1073</b>	-	-	-	-	-	-
IN	66	3	32	-	-	-	506	<b>750</b>	2	-
RB	448	79	31	2	1	<b>894</b>	-	521	-	-
RP	-	-	2	-	-	182	53	-	-	-
VBG	309	467	23	-	-	-	-	-	-	-
NNPS	22	1	521	-	-	-	-	-	-	-

Porcentaje de aciertos: 96,46 %

**Cuadro 6:** Segunda mitad de WSJ original vs segunda mitad WSJ etiquetado con TnT (entrenado con la primera mitad de WSJ)

TnT(1WSJ) WSJ	JJ	NN	VDN	VBD	NNP	RB	IN	NNPS	RP	VB
<b>JJ</b>	-	<b>1084</b>	610	79	564	352	44	2	18	83
<b>NN</b>	<b>1828</b>	-	36	30	<b>1102</b>	51	13	-	6	263
<b>VDN</b>	<b>829</b>	22	-	<b>874</b>	30	-	-	-	-	34
<b>VBD</b>	70	38	<b>1107</b>	-	13	-	-	-	-	21
<b>NNP</b>	454	458	18	5	-	40	19	<b>850</b>	2	18
<b>RB</b>	384	164	3	-	47	-	<b>753</b>	-	512	21
<b>IN</b>	43	5	-	-	19	<b>851</b>	-	-	<b>693</b>	4
<b>NNPS</b>	11	-	-	-	351	-	-	-	-	-
<b>RP</b>	2	1	-	-	1	85	227	-	-	-
<b>VB</b>	76	410	27	21	50	10	7	-	-	-

Porcentaje de aciertos: 96,20 %

**Cuadro 7:** Segunda mitad de WSJ original vs segunda mitad WSJ etiquetado con TnT (entrenado con primera mitad de WSJ)

TnT(1WSJ+NFI) WSJ	JJ	NN	NNP	VDN	VBD	IN	RB	RP	NNPS	VB
<b>JJ</b>	-	<b>842</b>	661	504	35	51	353	23	2	64
<b>NN</b>	<b>1808</b>	-	<b>1258</b>	34	30	10	40	4	-	217
<b>NNP</b>	334	318	-	14	8	11	27	1	<b>783</b>	18
<b>VDN</b>	<b>808</b>	21	26	-	<b>749</b>	-	-	-	-	26
<b>VBD</b>	46	16	13	<b>1072</b>	-	-	-	-	-	19
<b>IN</b>	43	3	23	-	-	-	652	<b>875</b>	-	3
<b>RB</b>	426	137	56	2	-	<b>878</b>	-	<b>692</b>	-	36
<b>RP</b>	2	1	-	-	-	188	45	-	-	-
<b>NNPS</b>	15	-	319	-	-	-	-	-	-	-
<b>VB</b>	58	272	37	24	29	7	8	-	-	-

Porcentaje de aciertos: 96,36 %

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 la primer y segunda mitad de WSJ etiquetado con TnT entrenado con la mitad restante

contra la misma mitad de WSJ etiquetado con TnT entrenado con la mitad restante + NFI.

**Cuadro 8:** *Primer mitad de WSJ etiquetado por TnT (entrenado con la segunda mitad WSJ) vs primer mitad WSJ etiquetado con TnT (entrenado con la segunda mitad de WSJ + NFI)*

TnT(2WSJ) \ TnT(2WSJ+NFI)	NN	JJ	VBN	VBD	NNP	VBG	VBP	VB	RP	RB
<b>NN</b>	-	<b>567</b>	32	19	<b>458</b>	<b>334</b>	108	<b>319</b>	-	43
<b>JJ</b>	<b>633</b>	-	143	39	<b>350</b>	63	12	59	2	162
<b>VBN</b>	15	<b>376</b>	-	<b>443</b>	7	-	5	8	-	-
<b>VBD</b>	10	36	<b>499</b>	-	13	-	5	8	-	-
<b>NNP</b>	159	142	1	2	-	6	3	18	-	5
<b>VBG</b>	175	216	-	1	36	-	-	1	-	2
<b>VBP</b>	94	2	-	7	2	-	-	233	-	3
<b>VB</b>	198	56	24	23	12	-	<b>326</b>	-	-	1
<b>RP</b>	-	-	-	-	-	-	-	-	-	6
<b>RB</b>	26	146	1	1	13	-	-	7	302	-

Porcentaje de aciertos: 98,39 %

**Cuadro 9:** *Segunda mitad WSJ etiquetado por TnT (entrenado con la primera mitad WSJ) vs segunda mitad WSJ etiquetado con TnT (entrenado con la primera mitad de WSJ + NFI)*

TnT(1WSJ) \ TnT(1WSJ+NFI)	JJ	NN	VBN	VBD	NNP	RP	RB	IN	VB	VBP
<b>JJ</b>	-	<b>510</b>	195	43	<b>319</b>	-	125	9	47	19
<b>NN</b>	<b>729</b>	-	27	25	<b>426</b>	1	41	4	<b>291</b>	69
<b>VBN</b>	<b>282</b>	20	-	<b>438</b>	4	-	1	-	11	5
<b>VBD</b>	56	12	<b>522</b>	-	3	-	1	-	8	3
<b>NNP</b>	117	133	5	3	-	-	10	9	25	4
<b>RP</b>	-	2	-	-	1	-	35	32	-	-
<b>RB</b>	160	35	-	-	35	<b>353</b>	-	<b>300</b>	16	2
<b>IN</b>	2	2	-	-	14	159	73	-	-	1
<b>VB</b>	43	183	18	17	17	1	4	1	-	275
<b>VBP</b>	16	102	3	7	1	-	1	1	201	-

Porcentaje de aciertos: 98,45 %

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 10:** *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	95.92 %
TnT entrenado con el primer 1/4 de WSJ + NFI	96.25 %
TnT entrenado con el segundo 1/4 de WSJ	95.88 %
TnT entrenado con el segundo 1/4 de WSJ + NFI	96.25 %
TnT entrenado con el tercer 1/4 de WSJ	95.90 %
TnT entrenado con el tercer 1/4 de WSJ + NFI	96.28 %
TnT entrenado con el cuarto 1/4 de WSJ	95.89 %
TnT entrenado con el cuarto 1/4 de WSJ + NFI	96.29 %

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 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:

- 95.31 % de acierto de etiquetas para el etiquetado de 9/10 de WSJ con TnT entrenado con 1/10 WSJ
- 96.09 % 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 corpus de entrenamiento que incorpora NFI.

### 0.1.3. Tercer experimento