



# Seminario de análisis del discurso

**Dr. César Antonio Aguilar**  
**Facultad de Lenguas y Letras**  
**20/09/2010**

**CAguilar@ingen.unam.mx**

# Estadísticas en corpus (1)

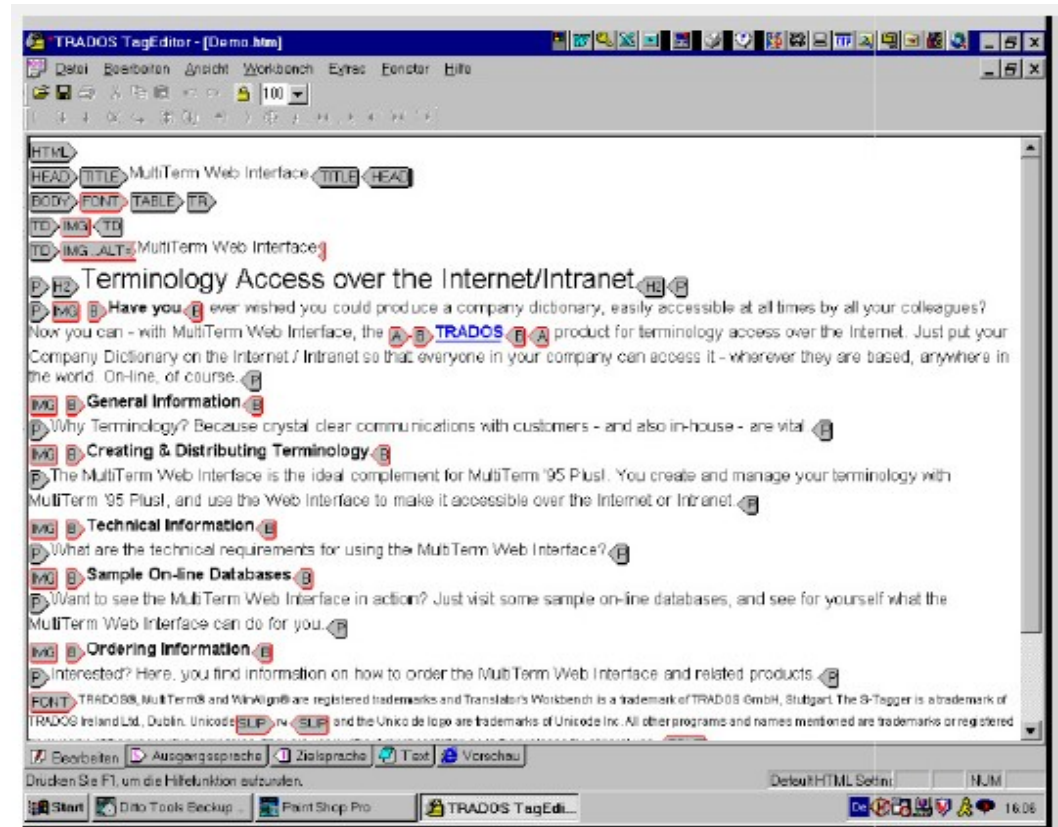
```
>>> reference = 'This is the reference data. Testing 123. aoacococ'
>>> test = 'Thos iz the rifirenci data. Testeng 123. aoaeoeoe'
>>> print ConfusionMatrix(reference, test)
| . 123T_acdefghinorstz |
+-----+
| <8> . . . . 1 . . . . . |
| . <2> . . . . . |
| 1 . <1> . . . . . |
| 2 . . <1> . . . . . |
| 3 . . . <1> . . . . . |
| T . . . . <2> . . . . . |
| _ . . . . . <4> . . . . . |
| a . . . . . <1> . . . . . |
| c . . . . . . <1> . . . . . |
| d . . . . . . . <1> . . . . . |
| e . . . . . . . . <6> . . 3 |
| f . . . . . . . . . <1> . . . . . |
| g . . . . . . . . . . <1> . . . . . |
| h . . . . . . . . . . . <2> . . . . . |
| i . . . . . 1 . . . <1> . 1 . . . . . |
| n . . . . . . . . . . . <2> . . . . . |
| o . . . . . . . . . . . . <3> . . . . . |
| r . . . . . . . . . . . . . <2> . . . . . |
| s . . . . . . . . . . . . . . <2> . 1 |
| t . . . . . . . . . . . . . . . <3> . . . . . |
| z . . . . . . . . . . . . . . . <4> . . . . . |
+-----+
(row = reference; col = test)
```

Como comentamos al inicio de esta unidad, uno de los rasgos interesantes del análisis de corpus es poder contar con un soporte empírico importante para validar alguna hipótesis sobre un fenómeno lingüístico.

# Estadísticas en corpus (2)

Con el ejercicio que realizaron dentro del corpus técnico del IULA, habrán visto lo pertinente que resulta usar algún método estadístico para calcular cuál es el grado de confianza que podemos esperar de un conjunto de datos para validar o no una hipótesis.

En terminología lo anterior es la regla: ¿cómo saber que una construcción sintáctica es un término? Simplemente calculando su frecuencia de uso dentro de documentos pertenecientes a un área técnica o científica.



# Estadísticas en corpus (3)

En estos casos, el concepto de **patrón** alude precisamente a una estructura lingüística que, debido a su regularidad estadística, es útil para encontrar información relevante para resolver una consulta.

	items	noise	Examples: acceptable term candidates	Examples: noise
<u>Individual queries:</u>				
P + N	2990	13,5%	<i>Hinterachse, Vollgas</i>	<i>Bedeutung, Gegensatz</i>
P + ADJ	1755	46%	<i>gefiltert, verstellbar</i>	<i>unabhängig, vorhanden</i>
P + V	1603	31%	<i>beschleunigen, einspritzen</i>	<i>beibehalten, entscheiden</i>
S + N	6349	8%	<i>Geschwindigkeit, Produktion</i>	<i>Kenntnis, Wirklichkeit</i>
S + ADJ	2674	29,5%	<i>bleifrei, lieferbar</i>	<i>gemeinsam, unmittelbar</i>
DC + N	5376	2%	<i>Dachluke, Motorbremse</i>	<i>Nachteil, Umgang</i>
DC + ADJ	786	22%	<i>thermodynamisch, betriebswarm</i>	<i>nachhaltig, nutzlos</i>
DC + V	585	29%	<i>einbauen, herausfiltern</i>	<i>verkraften, aufteilen</i>
Abbr.	783	1%	<i>km/h/s, FCKW-frei, TE-24</i>	<i>MANN, BOSCH-Ein-spritzpumpe</i>
<u>Combined Queries:</u>				
P + S + N	2270	13%	<i>Sonderausstattung, Verstellhebel</i>	<i>Verhältnis, Beschreibung</i>
P + S + ADJ	1508	51%	<i>einstellbar, hochporös</i>	<i>entsprechend, abhängig</i>
P + DC + N	1713	6,5%	<i>Leergas, Innenraumluft</i>	<i>Aufwand, Vorteil</i>
P + DC + ADJ	371	28%	<i>hochbelastet, ungefiltert</i>	<i>vorhanden, entsprechend</i>
P + DC + V	469	28%	<i>zuschalten, ansaugen</i>	<i>bewirken, erhalten</i>
S + DC + N	4233	2,5%	<i>Lasersensor, Kraftübertragung</i>	<i>Erteilung, Schimmelbildung</i>
S + DC + ADJ	747	20,5%	<i>elektronisch, nachgeschaltet</i>	<i>typisch, wirkungsvoll</i>
P + S + DC + N	1318	4%	<i>Abgasleitung, Feststellbremse</i>	<i>Erfahrung, Aufteilung</i>
P + S + DC + ADJ	363	24%	<i>vollelektronisch, hochbelastbar</i>	<i>ermittelt, vorausgegangen</i>

# Probabilidades (1)

Cuando manejamos enormes cantidades de datos lingüísticos (de cualquier tipo), si bien es cierto que muchas cosas las podemos deducir por sentido común, hay otras en donde requerimos determinar cuán probable o no es que tales datos reflejen un comportamiento determinado.

Un experimento sencillo: ¿cómo podemos calcular, en una partida de dados, cuál es el número más probable que salga en una tirada?

Delimitemos:

1. Un dado tiene 6 caras.
2. Cada cara representa un número del 1 al 6.
3. En teoría, en 6 tiradas me puede salir al menos un número del 1 al 6.



# Probabilidades (2)

Para hacer nuestros cálculos, necesitamos una tabla como la siguiente:

Tiradas	1	$1/6$	$1/6$	$1/6$	$1/6$	$1/6$	$1/6$
	2						
	3						
	4						
	5						
	6	1	2	3	4	5	6

Caras del dado



De primera instancia, podemos plantear dos cosas:

1. Si el dado no tiene ningún truco, en teoría cada tirada (en una serie de 6) nos muestra un número concreto del dado.
2. Por el contrario, si el dado tiene algún truco, es muy posible que tienda a mostrar un número una buena cantidad de veces.

## Probabilidades (3)

Dependiendo de una u otra opción, tenemos dos clases de probabilidades:

Probabilidad independiente o relativa: se da cuando planteamos la ocurrencia de dos eventos (a y B) que no mantienen ninguna relación entre sí.



$$P(A \text{ and } B) = P(A) \times P(B)$$

**Nota:** esto aplica para nuestro caso 1: cada vez que lancemos un dado, nos va a dar siempre una cara diferente.

# Probabilidades (4)

Probabilidad condicional: se da cuando planteamos la ocurrencia de dos eventos (A y B) o más, tomando en cuenta que el evento A determina el evento B.



Para decirlo en otras palabras: la probabilidad de que ocurra el evento B dado el evento A se puede representar con la fórmula:

$$p(B|A) = p(A \cap B) / p(A)$$

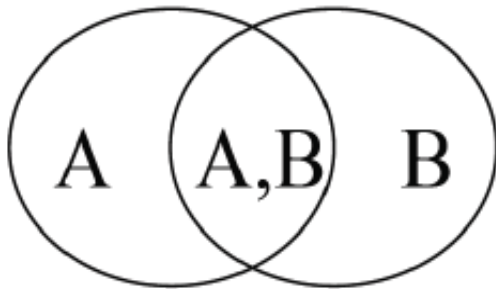
**Nota:** esto aplica para nuestro caso 2, y justo son la clase de probabilidades que nos interesan.



# Probabilidades (5)

Si nos ponemos más formales, una probabilidad condicional equivale a decir que:

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



Recordando las matemáticas de la secundaria, esto se puede representar con un diagrama de Venn:

Hay que notar que:

$$P(A, B) = P(A | B) \cdot P(B)$$

También:

$$P(A, B) = P(B, A)$$

# Probabilidades (6)

Veamos un ejemplo: con base en la probabilidad de que una figura sea una X, esto es,  $P(X)$ , ¿cuántas probabilidades hay de que dada una X sea roja, es decir,  $P(X|roja)$ ?

X	X	X	X	O	O
X	X	O	X	X	O
O	O	X	O	O	X
O	O	O	O	X	O
X	X	X	X	X	O

# Probabilidades y corpus (1)

Cuando vemos el lenguaje desde un punto de vista probabilístico, podemos considerar que los fenómenos que ocurren en él son de carácter condicional, más que fruto del azar.

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

Volviendo a algo que vimos la clase pasada: la detección de **n-gramas** es más que nada un proceso estadístico que ayuda al trabajo de inserción de etiquetas.

## Probabilidades y corpus (2)

Las probabilidades son empleadas muchas tareas, p. e., reconocer relaciones sintácticas y semánticas entre palabras:

En este caso, si establecemos que Shakespeare es un autor, es porque tal atributo aparece asociado a dicho nombre con una frecuencia alta, de acuerdo con la información que nos brinda, p. e., el **British National Corpus**.

“...works by such **authors** as Herrick, Goldsmith, and **Shakespeare**.”

“If you consider **authors** like **Shakespeare**...”

“Some **authors** (including **Shakespeare**)...”

“**Shakespeare** was the **author** of several...”

“**Shakespeare**, **author** of *The Tempest*...”



**Shakespeare** IS-A **author** (0.87)

# Experimento (1)

Veamos cómo podemos aprovechar esta clase de métodos para hacer búsquedas en corpus. En el siguiente sitio electrónico:

<http://labs.translated.net/>

En esta página encuentran un extractor de términos en línea, desarrollado por Translated Labs, que es una empresa dedicada a dar servicios de ingeniería lingüística.

Con este extractor, lo que queremos es obtener posibles candidatos a términos dentro de un documento, en este caso, un artículo escrito en inglés.

# Experimento (2)

El extractor nos muestra las 20 unidades que mostraron una alta frecuencia de uso que, por ello, pueden ser posibles candidatos a términos:

#	Extracted term	Score
1	<a href="#">verb</a>	62%
2	<a href="#">definitional patterns</a>	61%
3	<a href="#">definitional contexts</a>	60%
4	<a href="#">defining verb</a>	60%
5	<a href="#">como</a>	59%
6	<a href="#">cociente entre el área</a>	59%
7	<a href="#">con el intercambio entre</a>	59%
8	<a href="#">buque y el área</a>	59%
9	<a href="#">corpora</a>	58%
10	<a href="#">onomasiological dictionary</a>	57%
11	<a href="#">verbal predication</a>	57%
12	<a href="#">adverb como</a>	56%
13	<a href="#">con este trabajo</a>	55%
14	<a href="#">differentia</a>	55%
15	<a href="#">specialised corpora</a>	55%
16	<a href="#">definition aquel recurso</a>	55%
17	<a href="#">verbs</a>	55%
18	<a href="#">definir</a>	54%
19	<a href="#">defining verb lemmas</a>	54%
20	<a href="#">person plural pronoun</a>	53%

Podrán notar que no todos los posibles candidatos parecen ser buenos términos, pero algunos otros sí son buenas opciones.

Pregunta: en esta clase de tareas, ¿cómo podrían aprovechar probabilidades condicionales para explorar un corpus y obtener datos relevantes?

## Experimento (3)

Veamos otro caso, y tratemos de resolver estas preguntas:

1. En un artículo escrito por Yorick Wilks, ¿a qué se refieren los términos “semantic primitive derivation”, “statistical nlp” y “modern statistical nlp” ?
2. Los dos últimos términos, ¿se refieren al mismo concepto, o aluden a dos conceptos claramente diferenciados?
3. ¿Con qué palabras se relaciona el término “corpora”?
4. Las palabras “syntactic” y “computational”, ¿son buenos ejemplos de términos?

# Experimento (4)

Para tener una mejor idea, veamos los resultados:

#	Extracted term	Score
1	<a href="#">syntactic</a>	66%
2	<a href="#">linguistics</a>	63%
3	<a href="#">computational</a>	62%
4	<a href="#">corpora</a>	58%
5	<a href="#">wilks</a>	58%
6	<a href="#">sparck jones</a>	56%
7	<a href="#">syntactic structure</a>	56%
8	<a href="#">semantic primitive derivation</a>	55%
9	<a href="#">statistically-based empirical emphasis</a>	54%
10	<a href="#">programming language prolog</a>	54%
11	<a href="#">cat-dog sentence</a>	54%
12	<a href="#">syntactic analysers</a>	53%
13	<a href="#">dictionaries</a>	53%
14	<a href="#">computational linguistics</a>	53%
15	<a href="#">statistical nlp</a>	53%
16	<a href="#">grammar parsing</a>	53%
17	<a href="#">modern statistical nlp</a>	53%
18	<a href="#">overtax syntactic</a>	53%
19	<a href="#">computational semantics</a>	52%
20	<a href="#">roget 's thesaurus</a>	52%



# Tarea

Hagan una extracción de términos usando un artículo de Mira Ariel (el archivo se llama *Ariel01.txt*, y lo pueden descargar desde la página del blog. Una vez que obtengan sus términos, traten de responder a las siguientes preguntas:

1. Los terminos que enlista el extractor, ¿son todos los que hay en el artículo de Ariel? Revisen el artículo y propongan otros que no aparezcan.
2. ¿Cuál es el grado de confianza que le darían al extractor? Esto es: ¿los términos que no reconoció, son muchos, son pocos? ¿Cómo harían esta evaluación?
3. El término *pronoun* parece que es muy productivo en el texto de Ariel, y de hecho permite construir otros con una estructura más compleja (esto es, se constituyen a partir de 2 o más palabras). ¿Qué otros términos podrían derivar usando *pronoun*? Propongan por lo menos 10 posibles nuevos términos.

# Gracias por su atención

**Blog del curso:** <http://discurso-uaq.weebly.com/>