

Herramientas Libres

para

Procesamiento del Lenguaje Natural

Laura Alonso i Alemany

Facultad de Matemática, Astronomía y Física
UNC, Córdoba (Argentina)

<http://www.cs.famaf.unc.edu.ar/~laura>

5tas Jornadas Regionales de Software Libre
20 de noviembre de 2005

contenidos

qué es el PLN

aplicaciones

arquitectura

herramientas

preprocesos

análisis morfológico (*tagging*)

análisis sintáctico superficial (*chunking*)

análisis sintáctico (*parsing*)

análisis semántico

aplicaciones

cajas de herramientas

directorios de herramientas, recursos y documentación

aplicaciones de PLN

- ▶ recuperación de información
- ▶ acceso a bases de datos en lenguaje natural
- ▶ corrección automática (y sugerencia de palabras)
- ▶ resumen automático
- ▶ traducción automática
- ▶ soporte al aprendizaje de lenguas por computadora
- ▶ soporte a la descripción de lenguas por computadora
- ▶ creación semiautomática de mapas conceptuales
- ▶ detección de sentimientos

arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones
3. análisis morfológico
4. análisis sintáctico
5. análisis semántico

arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones

elgatomepescado

3. análisis morfológico
4. análisis sintáctico
5. análisis semántico

arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones

el gato come pescado

3. análisis morfológico
4. análisis sintáctico
5. análisis semántico

arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones
3. análisis morfológico
 - 3.1 detección de palabras especiales

Woody Allen llegó a Donosti el miércoles a las dos.

- 3.2 asignación de etiquetas
 - 3.3 desambiguación de etiquetas
4. análisis sintáctico
5. análisis semántico

arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones
3. análisis morfológico
 - 3.1 detección de palabras especiales

Woody Allen llegó a Donosti el miércoles a las dos.

- 3.2 asignación de etiquetas
 - 3.3 desambiguación de etiquetas
4. análisis sintáctico
5. análisis semántico

arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones
3. análisis morfológico

3.1 detección de palabras especiales

3.2 asignación de etiquetas

el	DA0MS0	el
gato	NCMS000	gato
come	VMIP3S0,VMPP2S0	comer
pescado	NCMS000,VMP00SM	pescado

3.3 desambiguación de etiquetas

4. análisis sintáctico
5. análisis semántico

arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones
3. análisis morfológico
 - 3.1 detección de palabras especiales
 - 3.2 asignación de etiquetas
 - 3.3 desambiguación de etiquetas

el	DA0MS0	el
gato	NCMS000	gato
come	VMIP3S0	comer
pescado	NCMS000	pescado

4. análisis sintáctico
5. análisis semántico

arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones
3. análisis morfológico
4. análisis sintáctico
 - 4.1 constituyentes básicos o *chunks*

el gato come pescado

- 4.2 estructura de oración
 - 4.3 funciones gramaticales, roles temáticos
5. análisis semántico

arquitectura básica de los sistemas de PLN

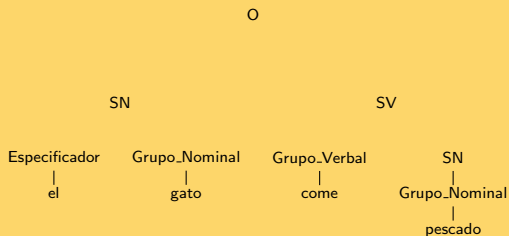
1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones
3. análisis morfológico
4. análisis sintáctico
 - 4.1 constituyentes básicos o *chunks*

Grupo_Nominal(el gato) Grupo_Verbal(come) Grupo_Nominal(pescado)

- 4.2 estructura de oración
 - 4.3 funciones gramaticales, roles temáticos
5. análisis semántico

arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones
3. análisis morfológico
4. análisis sintáctico
 - 4.1 constituyentes básicos o *chunks*
 - 4.2 estructura de oración



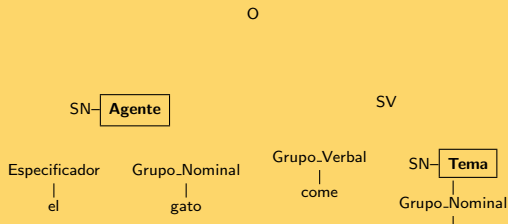
arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones
3. análisis morfológico
4. análisis sintáctico
 - 4.1 constituyentes básicos o *chunks*
 - 4.2 estructura de oración
 - 4.3 funciones gramaticales, roles temáticos



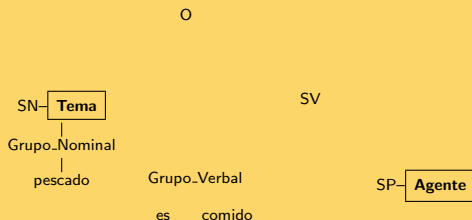
arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones
3. análisis morfológico
4. análisis sintáctico
 - 4.1 constituyentes básicos o *chunks*
 - 4.2 estructura de oración
 - 4.3 funciones gramaticales, roles temáticos



arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones
3. análisis morfológico
4. análisis sintáctico
 - 4.1 constituyentes básicos o *chunks*
 - 4.2 estructura de oración
 - 4.3 funciones gramaticales, roles temáticos



arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones
3. análisis morfológico
4. análisis sintáctico
5. análisis semántico

5.1 léxico

el gato	entidad → ser vivo → animal → ... → felino doméstico determinado
come	acción → voluntaria → ...
pescado	entidad → inanimado → natural → comestible entidad → ser vivo → animal → vertebrado → pez no determinado → masa

5.2 proposicional

arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones
3. análisis morfológico
4. análisis sintáctico
5. análisis semántico

5.1 léxico

Woody Allen	persona → artista → actor → cine
	persona → artista → director → cine
llegó	acción → desplazamiento → ...
a Donosti	lugar → ciudad
el miércoles a las dos	14:00GMT02/02/2005

5.2 proposicional

arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones
3. análisis morfológico
4. análisis sintáctico
5. análisis semántico

5.1 léxico

5.2 proposicional

$\exists \text{gato}(X) \wedge \exists \text{pescado}(Y) \wedge \text{come}(X,Y)$

arquitectura básica de los sistemas de PLN

1. reconocimiento de idioma
2. segmentación de palabras, oraciones, secciones
3. análisis morfológico
4. análisis sintáctico
5. análisis semántico
 - 5.1 léxico
 - 5.2 proposicional

llega(Woody_Allen,Donosti,14:00GMT02/02/2005)

arquitecturas simbólicas vs. probabilísticas

► arquitecturas simbólicas

1. un humano desarrolla reglas de análisis y/o diccionarios
2. el conocimiento codificado en las reglas y diccionarios se aplica mediante un analizador automático

► arquitecturas probabilísticas

1. uno (o más) humanos analizan una muestra representativa de lenguaje natural (*corpus anotado*)
2. se aplica un proceso de inferencia de conocimiento (reglas y/o diccionarios) a esta muestra
3. el conocimiento obtenido automáticamente en forma de reglas y diccionarios, a menudo probabilísticos, se aplica mediante un analizador automático

preprocesos para el análisis

- ▶ identificación de lengua
- ▶ segmentación de palabras (*tokenization*), oraciones, párrafos, secciones
- ▶ identificación de palabras especiales (fechas, distancias, etc.)
- ▶ identificación de entidades con nombre (*Named Entity Recognition*) (p.ej.: *Woody Allen*)
- ▶ lematización (o *stemming*)

el análisis morfológico

la mayor parte de herramientas asignan y desambiguan a la vez, y todas incluyen lematización

1. asignación de etiquetas
2. desambiguación de etiquetas

el análisis morfológico

la mayor parte de herramientas asignan y desambiguan a la vez, y todas incluyen lematización

1. asignación de etiquetas
2. desambiguación de etiquetas

el	DA0MS0	el
gato	NCMS000	gato
come	VMIP3S0,VMPP2S0	comer
pescado	NCMS000,VMP00SM	pescado

el análisis morfológico

la mayor parte de herramientas asignan y desambiguan a la vez, y todas incluyen lematización

1. asignación de etiquetas
2. desambiguación de etiquetas

el	DA0MS0	el
gato	NCMS000	gato
come	VMIP3S0	comer
pescado	NCMS000	pescado

diccionarios de palabras

- ▶ todos los analizadores morfológicos y sintácticos tienen un diccionario, en los casos de analizadores de código abierto, el diccionario es accesible
- ▶ para la lengua castellana, el diccionario de **Freeling** cubre un 90% de la lengua
- ▶ un importantísimo recurso léxico es **WordNet** y sus extensiones (**EuroWordNet**, **BalkaNet** y muchos otros), que veremos en la parte de análisis semántico.
- ▶ la **lista de lemarios del castellano** de Ismael Olea no tiene desperdicio.

correctores ortográficos

- ▶ **Xuxen** es un corrector ortográfico para el vasco
- ▶ **ispell** International Ispell is an interactive spell-checking program for Unix which supports a large number of European languages. An emacs interface is available as well as the standard command-line mode.
- ▶ **aspell** GNU Aspell is a Free and Open Source spell checker designed to eventually replace Ispell.
- ▶ el diccionario para el español **COES** está integrado en ispell y es de esperar que pronto lo esté en aspell.
- ▶ **myspell** es el corrector ortográfico de OpenOffice, basado en aspell.

qué es el PLN
herramientas

preprocesos
análisis morfológico (*tagging*)
análisis sintáctico superficial (*chunking*)
análisis sintáctico (*parsing*)
análisis semántico
aplicaciones
cajas de herramientas
directorios de herramientas, recursos y documentación

cómo se obtienen *taggers* probabilísticos

corpus de
entrenamiento

<i>el</i>	<i>gato</i>	<i>come</i>	<i>pescado</i>
<i>DA0MS0</i>	<i>NCMS000</i>	<i>VMIP3S0</i>	<i>NCMS000</i>

qué es el PLN
herramientas

preprocesos
análisis morfológico (*tagging*)
análisis sintáctico superficial (*chunking*)
análisis sintáctico (*parsing*)
análisis semántico
aplicaciones
cajas de herramientas
directorios de herramientas, recursos y documentación

cómo se obtienen *taggers* probabilísticos

**corpus de
entrenamiento**

**método de
inferencia**

*modelos ocultos de Markov (HMM),
modelos de máxima entropía, y otros*

qué es el PLN
herramientas

preprocesos
análisis morfológico (*tagging*)
análisis sintáctico superficial (*chunking*)
análisis sintáctico (*parsing*)
análisis semántico
aplicaciones
cajas de herramientas
directorios de herramientas, recursos y documentación

cómo se obtienen *taggers* probabilísticos

corpus de
entrenamiento

método de
inferencia

herramienta de
análisis

el-DA0MS0 gato-NCMS0 come VMIP3S0 VMPP2S0

qué es el PLN
herramientas

preprocesos
análisis morfológico (*tagging*)
análisis sintáctico superficial (*chunking*)
análisis sintáctico (*parsing*)
análisis semántico
aplicaciones
cajas de herramientas
directorios de herramientas, recursos y documentación

cómo se obtienen *taggers* probabilísticos

corpus de
entrenamiento

método de
inferencia

herramienta de
análisis

el-DA0MS0 gato-NCMS0 **come VMIP3S0 VMPP2S0**
—analizador→ come-VMIP3S0

corpus anotados

para el español:

- ▶ **3Ib** un corpus de 100.000 palabras anotadas con su categoría morfosintáctica y su interpretación sintáctica. También hay corpus de 50.000 palabras del catalán y del euskera. Libre para investigación.

para otras lenguas:

- ▶ **Susanne** es un extracto de 130.000 palabras del corpus Brown de inglés americano, analizadas sintácticamente
- ▶ **Christine** es un extracto de 80.000 palabras de lenguaje oral del corpus del inglés British National Corpus, analizadas sintácticamente
- ▶ **Lucy** es un corpus del inglés británico de 165.000 palabras, analizadas sintácticamente

analizadores morfológicos

- ▶ **Stanford POS tagger** java, código abierto (GPL). Se incluyen dos modelos para el inglés.
- ▶ **Brill's Transformation-based learning Tagger**
- ▶ **Maximum Entropy part of speech tagger MXPOST**
- ▶ **TnT**
- ▶ **SVMTool**
- ▶ **Original Xerox Tagger**
- ▶ **ACOPOST**
- ▶ **μ -TBL**
- ▶ **QTA**
- ▶ **The TOSCA/LOB tagger**
- ▶ **TreeTagger**
- ▶ **Lingua-EN-Tagger**

analizadores morfológicos

- ▶ Stanford POS tagger
- ▶ Brill's Transformation-based learning Tagger C, aproximación simbólica.
- ▶ Maximum Entropy part of speech tagger MXPOST
- ▶ TnT
- ▶ SVMTool
- ▶ Original Xerox Tagger
- ▶ ACOPOST
- ▶ μ -TBL
- ▶ QTA
- ▶ The TOSCA/LOB tagger
- ▶ TreeTagger
- ▶ Lingua-EN-Tagger

analizadores morfológicos

- ▶ Stanford POS tagger
- ▶ Brill's Transformation-based learning Tagger
- ▶ Maximum Entropy part of speech tagger MXPOST java
(Archivos de clases, no fuente). Incluye un detector de finales de oración.
- ▶ TnT
- ▶ SVMTool
- ▶ Original Xerox Tagger
- ▶ ACOPOST
- ▶ μ -TBL
- ▶ QTA
- ▶ The TOSCA/LOB tagger
- ▶ TreeTagger

analizadores morfológicos

- ▶ Stanford POS tagger
- ▶ Brill's Transformation-based learning Tagger
- ▶ Maximum Entropy part of speech tagger MXPOST
- ▶ TnT para Solaris y Linux. Muy eficiente. Incluye modelos para inglés y alemán. Licencia de uso libre para fines no comerciales.
- ▶ SVMTool
- ▶ Original Xerox Tagger
- ▶ ACOPOST
- ▶ μ -TBL
- ▶ QTA
- ▶ The TOSCA/LOB tagger
- ▶ TreeTagger

analizadores morfológicos

- ▶ Stanford POS tagger
- ▶ Brill's Transformation-based learning Tagger
- ▶ Maximum Entropy part of speech tagger MXPOST
- ▶ TnT
- ▶ SVMTool C y Perl, código abierto (LGPL). Se basa en support vector machines, incorpora modelos para español, catalán e inglés.
- ▶ Original Xerox Tagger
- ▶ ACOPOST
- ▶ μ -TBL
- ▶ QTA
- ▶ The TOSCA/LOB tagger
- ▶ TreeTagger

analizadores morfológicos

- ▶ Stanford POS tagger
- ▶ Brill's Transformation-based learning Tagger
- ▶ Maximum Entropy part of speech tagger MXPOST
- ▶ TnT
- ▶ SVMTool
- ▶ Original Xerox Tagger Common lisp, basado en HMM.
También hay una **versión para el español**.
- ▶ ACOPOST
- ▶ μ -TBL
- ▶ QTA
- ▶ The TOSCA/LOB tagger
- ▶ TreeTagger
- ▶ Lingua-EN-Tagger

analizadores morfológicos

- ▶ Stanford POS tagger
- ▶ Brill's Transformation-based learning Tagger
- ▶ Maximum Entropy part of speech tagger MXPOST
- ▶ TnT
- ▶ SVMTool
- ▶ Original Xerox Tagger
- ▶ ACOPOST varios analizadores para POSIX, en C y Perl, código abierto (GNU). Implementa modelos de máxima entropía, de aprendizaje basado en transformaciones y de HMM de 3 palabras.
- ▶ μ -TBL
- ▶ QTA
- ▶ The TOSCA/LOB tagger

analizadores morfológicos

- ▶ Stanford POS tagger
- ▶ Brill's Transformation-based learning Tagger
- ▶ Maximum Entropy part of speech tagger MXPOST
- ▶ TnT
- ▶ SVMTool
- ▶ Original Xerox Tagger
- ▶ ACOPOST
- ▶ μ -TBL Prolog, aprendizaje basado en transformaciones, también se puede usar para otro tipo de aprendizaje.
- ▶ QTA
- ▶ The TOSCA/LOB tagger
- ▶ TreeTagger
- ▶ Lingua-EN-Tagger

analizadores morfológicos

- ▶ Stanford POS tagger
- ▶ Brill's Transformation-based learning Tagger
- ▶ Maximum Entropy part of speech tagger MXPOST
- ▶ TnT
- ▶ SVMTool
- ▶ Original Xerox Tagger
- ▶ ACOPOST
- ▶ μ -TBL
- ▶ QTA java (Archivos de clases, no fuente). Basado en HMM. Incluye diccionarios del inglés y del alemán.
- ▶ The TOSCA/LOB tagger
- ▶ TreeTagger
- ▶ Lingua-EN-Tagger

analizadores morfológicos

- ▶ Stanford POS tagger
- ▶ Brill's Transformation-based learning Tagger
- ▶ Maximum Entropy part of speech tagger MXPOST
- ▶ TnT
- ▶ SVMTool
- ▶ Original Xerox Tagger
- ▶ ACOPOST
- ▶ μ -TBL
- ▶ QTA
- ▶ The TOSCA/LOB tagger sistema histórico, sólo para MS-DOS.
- ▶ TreeTagger
- ▶ Lingua-EN-Tagger

analizadores morfológicos

- ▶ Stanford POS tagger
- ▶ Brill's Transformation-based learning Tagger
- ▶ Maximum Entropy part of speech tagger MXPOST
- ▶ TnT
- ▶ SVMTool
- ▶ Original Xerox Tagger
- ▶ ACOPOST
- ▶ μ -TBL
- ▶ QTA
- ▶ The TOSCA/LOB tagger
- ▶ TreeTagger con diccionarios para inglés, alemán, francés e italiano. Para Solaris y Linux. Basado en árboles de decisión.
- ▶ Lingua-EN-Tagger

analizadores morfológicos

- ▶ Stanford POS tagger
- ▶ Brill's Transformation-based learning Tagger
- ▶ Maximum Entropy part of speech tagger MXPOST
- ▶ TnT
- ▶ SVMTool
- ▶ Original Xerox Tagger
- ▶ ACOPOST
- ▶ μ -TBL
- ▶ QTA
- ▶ The TOSCA/LOB tagger
- ▶ TreeTagger
- ▶ Lingua-EN-Tagger Perl, basado en HMM de 2 palabras.
- ▶ PoSTech Korean morphological analyzer and tagger

análisis de grupos lingüísticos

se identifican grupos lingüísticos o *chunks*: p.ej., [*el gato*] [*come*] [*pescado*]:

- ▶ **FreeLing** código abierto (LGPL), con diccionarios y gramáticas para español, catalán e inglés. Tiene un diccionario del español que cubre más del 90% de la lengua, el diccionario de más cobertura de uso totalmente libre.
- ▶ **YamCha**
- ▶ **LingPipe**
- ▶ **fnTBL**

análisis de grupos lingüísticos

se identifican grupos lingüísticos o *chunks*: p.ej., *[el gato] [come] [pescado]*:

- ▶ **FreeLing**
- ▶ **YamCha** C/C++ código abierto, para el inglés, ganador de un concurso en reconocimiento de entidades con nombre (p.ej.: *Woody Allen*)
- ▶ **LingPipe**
- ▶ **fnTBL**

análisis de grupos lingüísticos

se identifican grupos lingüísticos o *chunks*: p.ej., *[el gato] [come] [pescado]*:

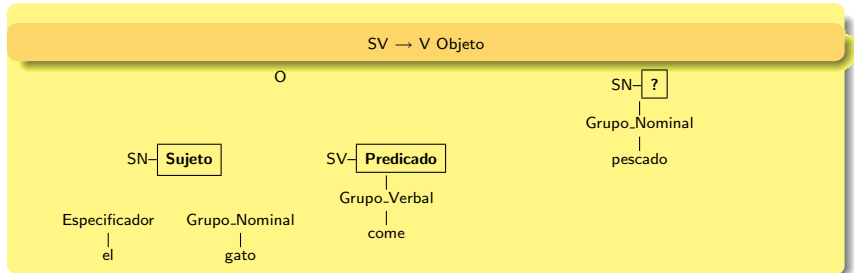
- ▶ **FreeLing**
- ▶ **YamCha**
- ▶ **LingPipe** java (GPL), reconoce entidades con nombre, finales de oración, e incluso co-referencia dentro de un documento
- ▶ **fnTBL**

análisis sintáctico tradicional (manual)

1. uno (o más) lingüistas crean una gramática de la lengua
 - ▶ reglas independientes de contexto (*Context Free Grammar*)
 $SN \rightarrow Det N$
 - ▶ reglas enriquecidas con rasgos (*Unification Grammar*)
 $SN_{fem,sg} \rightarrow Det_{fem,sg} N_{fem,sg}$
 - ▶ basada en el léxico (*Lexicalized Grammar*)
 $SN_{gato} \rightarrow Det N_{gato}$
2. un analizador (o *parser*) utiliza esta gramática para asignar estructura a oraciones no vistas previamente

análisis sintáctico tradicional (manual)

1. uno (o más) lingüistas crean una gramática de la lengua
2. un analizador (o *parser*) utiliza esta gramática para asignar estructura a oraciones no vistas previamente



análisis sintáctico basado en gramáticas manuales

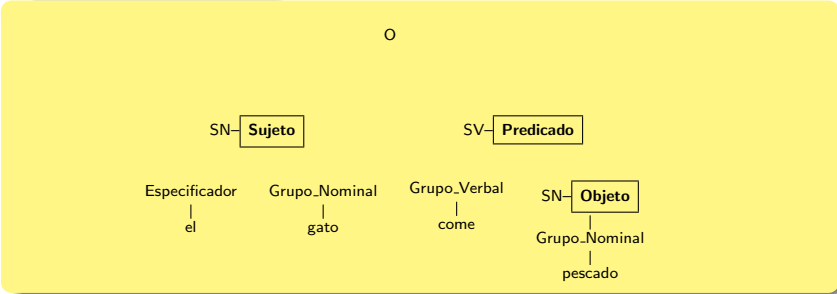
- ▶ Prolog tiene una extensión para implementar gramáticas libres de contexto: DCG (Definite Clause Grammars)
- ▶ **ALE** es un analizador para gramáticas de unificación basada en prolog, incluye gramáticas del inglés en HPSG (una clase famosa de gramáticas de unificación)
- ▶ **Link Grammar C**, código abierto, basada en formalismo de dependencias
- ▶ **English Resource Grammar** gramática HPSG del inglés, funciona sobre **LKB**
- ▶ **Jacy** gramática HPSG del japonés
- ▶ **Modern Greek Resource Grammar** gramática HPSG para el griego moderno

qué es el PLN
herramientas

- preprocesos
- análisis morfológico (*tagging*)
- análisis sintáctico superficial (*chunking*)
- análisis sintáctico (*parsing*)**
- análisis semántico
- aplicaciones
- cajas de herramientas
- directorios de herramientas, recursos y documentación

cómo se obtienen *parsers* probabilísticos

**corpus de
entrenamiento**



qué es el PLN
herramientas

preprocesos
análisis morfológico (*tagging*)
análisis sintáctico superficial (*chunking*)
análisis sintáctico (*parsing*)
análisis semántico
aplicaciones
cajas de herramientas
directorios de herramientas, recursos y documentación

cómo se obtienen *parsers* probabilísticos

**corpus de
entrenamiento**

**método de
inferencia**

gramáticas libres de contexto probabilísticas (lexicalizadas)
(Probabilistic (lexicalized) Context Free Grammars)

qué es el PLN
herramientas

preprocesos
análisis morfológico (*tagging*)
análisis sintáctico superficial (*chunking*)
análisis sintáctico (*parsing*)
análisis semántico
aplicaciones
cajas de herramientas
directorios de herramientas, recursos y documentación

cómo se obtienen *parsers* probabilísticos

corpus de
entrenamiento

método de
inferencia

herramienta de
análisis

$SV \rightarrow V \text{ Objeto } P = .82$
 $SV \rightarrow V \text{ Circunstancial } P = .18$



analizadores probabilísticos: corpus anotados

la mayor parte de corpus son pagos, excepto unos pocos, que son chicos :(

- ▶ 3Ib
- ▶ Susanne
- ▶ Christine
- ▶ Lucy

analizadores sintácticos probabilísticos

- ▶ software by Mark Johnson
- ▶ MINIPAR
- ▶ Stanford Lexicalized Parser
- ▶ Eugene Charniak's parser
- ▶ Michael Collins' parser
- ▶ Dan Bikel's parser
- ▶ Apple Pie Parser

analizadores sintácticos probabilísticos

- ▶ **software by Mark Johnson** Mark Johnson tiene disponible en su página web un montón de software relacionado con parsing, incluyendo un parser basado en reranking del 2005, una implementación en C muy eficiente de un parser clásico (CKY) y un parser muy popular en common lisp
- ▶ **MINIPAR**
- ▶ **Stanford Lexicalized Parser**
- ▶ **Eugene Charniak's parser**
- ▶ **Michael Collins' parser**
- ▶ **Dan Bikel's parser**
- ▶ **Apple Pie Parser**

analizadores sintácticos probabilísticos

- ▶ software by Mark Johnson
- ▶ MINIPAR C++, código abierto, para el inglés, muy eficiente y muy claro
- ▶ Stanford Lexicalized Parser
- ▶ Eugene Charniak's parser
- ▶ Michael Collins' parser
- ▶ Dan Bikel's parser
- ▶ Apple Pie Parser

analizadores sintácticos probabilísticos

- ▶ software by Mark Johnson
- ▶ MINIPAR
- ▶ Stanford Lexicalized Parser java, código abierto, para el inglés
- ▶ Eugene Charniak's parser
- ▶ Michael Collins' parser
- ▶ Dan Bikel's parser
- ▶ Apple Pie Parser

analizadores sintácticos probabilísticos

- ▶ software by Mark Johnson
- ▶ MINIPAR
- ▶ Stanford Lexicalized Parser
- ▶ Eugene Charniak's parser C++, código abierto, para el inglés
- ▶ Michael Collins' parser
- ▶ Dan Bikel's parser
- ▶ Apple Pie Parser

analizadores sintácticos probabilísticos

- ▶ software by Mark Johnson
- ▶ MINIPAR
- ▶ Stanford Lexicalized Parser
- ▶ Eugene Charniak's parser
- ▶ Michael Collins' parser C, fuente y ejecutables, para el inglés, también existe una versión que se puede correr como un daemon, documentación de su adaptación al checo
- ▶ Dan Bikel's parser
- ▶ Apple Pie Parser

analizadores sintácticos probabilísticos

- ▶ software by Mark Johnson
- ▶ MINIPAR
- ▶ Stanford Lexicalized Parser
- ▶ Eugene Charniak's parser
- ▶ Michael Collins' parser
- ▶ Dan Bikel's parser java, código abierto y clases, incluye una reimplementación exacta del parser de Collins y packs para inglés, chino y árabe, e está trabajando en adaptaciones al español y al coreano
- ▶ Apple Pie Parser

qué entendemos por análisis semántico?

Woody Allen	persona → artista → actor → cine
	persona → artista → director → cine
llegó	acción → desplazamiento → ...
a Donosti	lugar → ciudad
el miércoles a las dos	14:00GMT02/02/2005

para ello hay que asociar cada palabra a un *sentido*

diccionarios de sentidos y ontologías

Existen diversos diccionarios de sentidos organizados en forma de árbol (*ontologías léxicas*):

- ▶ WordNet
- ▶ EuroWordNet

diccionarios de sentidos y ontologías

Existen diversos diccionarios de sentidos organizados en forma de árbol (*ontologías léxicas*):

- ▶ **WordNet** 155.00 nombres, verbos y adjetivos del inglés se organizan en grupos de sinónimos (*synsets*) que a su vez se relacionan entre ellos mediante relaciones semánticas: *tipo de*, *contrario de*, etc. Totalmente libre, en varios formatos de uso y consulta y con extensa documentación, científica y técnica.
- ▶ **EuroWordNet**

diccionarios de sentidos y ontologías

Existen diversos diccionarios de sentidos organizados en forma de árbol (*ontologías léxicas*):

- ▶ **WordNet**
- ▶ **EuroWordNet** usando la estructura de WordNet como esqueleto común (Inter-Lingual-Index, ILI) se construyen ontologías para español, holandés, italiano, alemán, francés, checo y estonio. Libres para uso no comercial

diccionarios de sentidos y ontologías

Existen diversos diccionarios de sentidos organizados en forma de árbol (*ontologías léxicas*):

- ▶ WordNet
- ▶ EuroWordNet

y también **algoritmos para la asignación de palabras a sentidos** basados en WordNet

recuperación de información (*information retrieval*)

- ▶ Search Tools
- ▶ IN TeraScale Retrieval
- ▶ REtrieval COmponent INtegrator
- ▶ The Lemur Toolkit

recuperación de información (*information retrieval*)

- ▶ **Search Tools** un directorio que ayuda a encontrar el motor de búsqueda (*search engine*) más adecuado para cada necesidad: para web, intranets, diferentes tipos de datos, de aplicación, etc., con un apartado especial para **motores de código abierto**, incluyendo un artículo comparativo.
- ▶ **IN TeraScale Retrieval**
- ▶ **REtrieval COmponent INtegrator**
- ▶ **The Lemur Toolkit**

recuperación de información (*information retrieval*)

- ▶ Search Tools
- ▶ IN TeraScale Retrieval C++, GNU, un toolkit completo de herramientas de IR para todos los sistemas POSIX, con énfasis en recuperación de información semiestructurada (HTML, XML)
- ▶ REtrieval COmponent INtegrator
- ▶ The Lemur Toolkit

recuperación de información (*information retrieval*)

- ▶ Search Tools
- ▶ IN TeraScale Retrieval
- ▶ REtrieval COmponent INtegrator herramientas libres para desarrollar investigación en recuperación de información
- ▶ The Lemur Toolkit

recuperación de información (*information retrieval*)

- ▶ Search Tools
- ▶ IN TeraScale Retrieval
- ▶ REtrieval COmponent INtegrator
- ▶ The Lemur Toolkit explota el trabajo de modelado de lenguaje en otras áreas de PLN para aplicarlo a recuperación de información, orientado sobretodo a investigación

traducción automática (*machine translation*)

- ▶ Apertium
- ▶ Delph-In
- ▶ Laurie's links
- ▶ la serie de workshops sobre Teaching Machine Translation (con interesantes artículos sobre recursos libres): 2001, 2003

traducción automática (*machine translation*)

- ▶ **Apertium** un traductor entre lenguas romances de España, código abierto, basado en análisis superficial dentro del proyecto **OpenTrad**, que también desarrolla un traductor de código abierto entre castellano y euskera, basado en análisis sintáctico completo
- ▶ **Delph-In**
- ▶ **Laurie's links**
- ▶ la serie de workshops sobre Teaching Machine Translation (con interesantes artículos sobre recursos libres): **2001**, **2003**

traducción automática (*machine translation*)

- ▶ **Apertium**
- ▶ **Delph-In** es un proyecto de comprensión profunda de lenguaje natural cuyos recursos (libres!) han sido aplicados a traducción automática
- ▶ **Laurie's links**
- ▶ la serie de workshops sobre Teaching Machine Translation (con interesantes artículos sobre recursos libres): **2001**, **2003**

traducción automática (*machine translation*)

- ▶ Apertium
- ▶ Delph-In
- ▶ Laurie's links una exhaustiva lista de motores de traducción disponibles via web (en el año 2000), donde se especifica los idiomas que tratan, el texto máximo permitido, etc.
- ▶ la serie de workshops sobre Teaching Machine Translation (con interesantes artículos sobre recursos libres): 2001, 2003

traducción automática estadística y corpus alineados

La idea básica de los sistemas de traducción automática estadística es obtener un diccionario bilingüe a partir de corpus paralelos en las dos lenguas, que han sido alineados.

traducción automática estadística y corpus alineados

La idea básica de los sistemas de traducción automática estadística es obtener un diccionario bilingüe a partir de corpus paralelos en las dos lenguas, que han sido alineados.

<i>el</i>	<i>the</i>
<i>gato</i>	<i>cat</i>
<i>come</i>	<i>eats</i>
<i>pescado</i>	<i>fish</i>

traducción automática estadística y corpus alineados

La idea básica de los sistemas de traducción automática estadística es obtener un diccionario bilingüe a partir de corpus paralelos en las dos lenguas, que han sido alineados.

- ▶ **Europarl** corpus de documentos de la Unión Europea, con cerca de 20 millones de palabras en total, con unas 740.000 frases de cada una de las 11 lenguas, alineados manualmente a nivel de oración.
- ▶ **Hansards**
- ▶ **CRATER**
- ▶ **OPUS**
- ▶ **GNOME's GUI messages translation statistics**
- ▶ **Emille**

traducción automática estadística y corpus alineados

La idea básica de los sistemas de traducción automática estadística es obtener un diccionario bilingüe a partir de corpus paralelos en las dos lenguas, que han sido alineados.

- ▶ **Europarl**
- ▶ **Hansards** corpus de documentos del parlamento de Canadá, paralelos en inglés y francés, alineados a nivel de oración o menor
- ▶ **CRATER**
- ▶ **OPUS**
- ▶ **GNOME's GUI messages translation statistics**
- ▶ **Emille**

traducción automática estadística y corpus alineados

La idea básica de los sistemas de traducción automática estadística es obtener un diccionario bilingüe a partir de corpus paralelos en las dos lenguas, que han sido alineados.

- ▶ **Europarl**
- ▶ **Hansards**
- ▶ **CRATER** corpus alineado trilingüe: inglés, francés y castellano, con anotaciones morfosintácticas revisadas manualmente.
- ▶ **OPUS**
- ▶ **GNOME's GUI messages translation statistics**
- ▶ **Emille**

traducción automática estadística y corpus alineados

La idea básica de los sistemas de traducción automática estadística es obtener un diccionario bilingüe a partir de corpus paralelos en las dos lenguas, que han sido alineados.

- ▶ **Europarl**
- ▶ **Hansards**
- ▶ **CRATER**
- ▶ **OPUS** corpus de manuales técnicos (e.g., software libre, documentación de la Unión Europea) traducidos a varios idiomas, que han sido alineados automáticamente, están disponibles el corpus y el software de alineación.
- ▶ **GNOME's GUI messages translation statistics**
- ▶ **Emille**

traducción automática estadística y corpus alineados

La idea básica de los sistemas de traducción automática estadística es obtener un diccionario bilingüe a partir de corpus paralelos en las dos lenguas, que han sido alineados.

- ▶ **Europarl**
- ▶ **Hansards**
- ▶ **CRATER**
- ▶ **OPUS**
- ▶ **GNOME's GUI messages translation statistics**
- ▶ **Emille** corpus de 200.000 palabras, paralelo en inglés, hindi, bengalí, punjabí, gujarati y urdu.

traducción automática estadística y corpus alineados

La idea básica de los sistemas de traducción automática estadística es obtener un diccionario bilingüe a partir de corpus paralelos en las dos lenguas, que han sido alineados.

- ▶ **EGYPT** es un toolkit para desarrollar sistemas de traducción automática estadística a partir de corpus paralelos.
- ▶ **Rada Mihalcea** mantiene un extenso repositorio sobre alineación para traducción automática: corpus, software, etc.

exploración de datos: clasificación y clustering

- ▶ **R** la versión libre de S: un software para computación estadística y gráficos. Para todas las plataformas.
- ▶ **The 'Bow' Toolkit** librerías en C para análisis estadístico de texto, modelado de lenguaje y recuperación de información.
- ▶ **Weka** algoritmos para todo tipo de tareas de minería de datos, se pueden usar desde terminal, desde interfaz gráfica o desde tu propio código java. Cuenta con un libro de introducción a Weka y a la minería de datos en general y una activa lista de usuarios.
- ▶ **Mallet** es una herramienta para aplicar todo tipo de técnicas de Machine Learning a lenguaje natural

entornos para crear analizadores

- ▶ NLTK - Natural Language Toolkit
- ▶ GATE - a General Architecture for Text Engineering
- ▶ CCG Library
- ▶ EDG
- ▶ LKB
- ▶ Matrix
- ▶ NLPfarm
- ▶ Ellogon

entornos para crear analizadores

- ▶ **NLTK - Natural Language Toolkit** una suite de librerías y programas en Python para desarrollar gramáticas y analizadores de todo tipo, simbólico y estadístico
- ▶ **GATE - a General Architecture for Text Engineering**
- ▶ **CCG Library**
- ▶ **EDG**
- ▶ **LKB**
- ▶ **Matrix**
- ▶ **NLPfarm**
- ▶ **Ellogon**

entornos para crear analizadores

- ▶ **NLTK - Natural Language Toolkit**
- ▶ **GATE - a General Architecture for Text Engineering** java, código abierto, muy bien documentado, resultado de un gran proyecto, ampliamente usado para diversas tareas de PLN, sobretodo orientado a comprensión profunda
- ▶ **CCG Library**
- ▶ **EDG**
- ▶ **LKB**
- ▶ **Matrix**
- ▶ **NLPfarm**
- ▶ **Ellogon**

entornos para crear analizadores

- ▶ **NLTK - Natural Language Toolkit**
- ▶ **GATE - a General Architecture for Text Engineering**
- ▶ **CCG Library** una colección de herramientas para desarrollar analizadores en el marco de Combinatory Categorical Grammar, java, código abierto, LGPL o librería GNU
- ▶ **EDG**
- ▶ **LKB**
- ▶ **Matrix**
- ▶ **NLPfarm**
- ▶ **Ellogon**

entornos para crear analizadores

- ▶ NLTK - Natural Language Toolkit
- ▶ GATE - a General Architecture for Text Engineering
- ▶ CCG Library
- ▶ EDG Example-based Development of Grammars, un sistema en lisp para desarrollar analizadores en el marco de Head Driven Phrase Structure Grammar (HPSG)
- ▶ LKB
- ▶ Matrix
- ▶ NLPfarm
- ▶ Ellogon

entornos para crear analizadores

- ▶ NLTK - Natural Language Toolkit
- ▶ GATE - a General Architecture for Text Engineering
- ▶ CCG Library
- ▶ EDG
- ▶ LKB entorno para desarrollar gramáticas y léxicos basados en gramáticas de unificación, explotando los principios de estructuras tipadas del proyecto DELPH-IN
- ▶ Matrix
- ▶ NLPfarm
- ▶ Ellogon

entornos para crear analizadores

- ▶ NLTK - Natural Language Toolkit
- ▶ GATE - a General Architecture for Text Engineering
- ▶ CCG Library
- ▶ EDG
- ▶ LKB
- ▶ **Matrix** un kit de principiante para desarrollar gramáticas HPSG en LKB
- ▶ NLPfarm
- ▶ Ellogon

entornos para crear analizadores

- ▶ NLTK - Natural Language Toolkit
- ▶ GATE - a General Architecture for Text Engineering
- ▶ CCG Library
- ▶ EDG
- ▶ LKB
- ▶ Matrix
- ▶ NLPfarm concentra diversos módulos para procesamiento de diálogo en java
- ▶ Ellogon

entornos para crear analizadores

- ▶ NLTK - Natural Language Toolkit
- ▶ GATE - a General Architecture for Text Engineering
- ▶ CCG Library
- ▶ EDG
- ▶ LKB
- ▶ Matrix
- ▶ NLPfarm
- ▶ Ellogon entorno gráfico multiplataforma para todo tipo de aplicaciones de ingeniería del lenguaje

entornos para crear corpus anotados

- ▶ **Alembic** un banco de trabajo (*workbench*) para desarrollar corpus anotados y analizadores que se basen en ellos con una gran reducción del esfuerzo humano
- ▶ **Wordfreak** una herramienta de anotación java (mozilla public license 1.1), para anotaciones humanas, automáticas y semiautomáticas (mediante active learning)
- ▶ **AGTK** herramienta para anotar señales acústicas y todo tipo de series temporales (audio, video), basada en grafos

directorios de herramientas, recursos y documentación

- ▶ OpenNLP
- ▶ el grupo de PLN de Stanford
- ▶ Kenji Kita
- ▶ Manuel Barberá
- ▶ recursos del Summer Institute of Linguistics
- ▶ recursos de la Linguist List
- ▶ WEBIR/IE

directorios de herramientas, recursos y documentación

- ▶ **OpenNLP** es un directorio de recursos de PLN de código abierto en sourceforge
- ▶ el grupo de PLN de Stanford
- ▶ Kenji Kita
- ▶ Manuel Barberá
- ▶ recursos del Summer Institute of Linguistics
- ▶ recursos de la Linguist List
- ▶ WEBIR/IE

directorios de herramientas, recursos y documentación

- ▶ OpenNLP
- ▶ el grupo de PLN de Stanford mantiene lista de recursos y herramientas de PLN probabilísticas y de lingüística computacional muy extenso y actualizado
- ▶ Kenji Kita
- ▶ Manuel Barberá
- ▶ recursos del Summer Institute of Linguistics
- ▶ recursos de la Linguist List
- ▶ WEBIR/IE

directorios de herramientas, recursos y documentación

- ▶ OpenNLP
- ▶ el grupo de PLN de Stanford
- ▶ Kenji Kita también tiene una extensa página de links a recursos y herramientas para PLN
- ▶ Manuel Barberá
- ▶ recursos del Summer Institute of Linguistics
- ▶ recursos de la Linguist List
- ▶ WEBIR/IE

directorios de herramientas, recursos y documentación

- ▶ OpenNLP
- ▶ el grupo de PLN de Stanford
- ▶ Kenji Kita
- ▶ Manuel Barberá también mantiene una muy respetable lista de enlaces, poco actualizados pero con el interés de centrarse bastante en lenguas romances
- ▶ recursos del Summer Institute of Linguistics
- ▶ recursos de la Linguist List
- ▶ WEBIR/IE

directorios de herramientas, recursos y documentación

- ▶ OpenNLP
- ▶ el grupo de PLN de Stanford
- ▶ Kenji Kita
- ▶ Manuel Barberá
- ▶ recursos del Summer Institute of Linguistics orientados sobretodo a la descripción de lenguas
- ▶ recursos de la Linguist List
- ▶ WEBIR/IE

directorios de herramientas, recursos y documentación

- ▶ OpenNLP
- ▶ el grupo de PLN de Stanford
- ▶ Kenji Kita
- ▶ Manuel Barberá
- ▶ recursos del Summer Institute of Linguistics
- ▶ recursos de la [Linguist List](#) cubren todo el espectro de la lingüística: descripción, aprendizaje, diccionarios, fonética, lingüística histórica... y por supuesto PLN
- ▶ [WEBIR/IE](#)

directorios de herramientas, recursos y documentación

- ▶ OpenNLP
- ▶ el grupo de PLN de Stanford
- ▶ Kenji Kita
- ▶ Manuel Barberá
- ▶ recursos del Summer Institute of Linguistics
- ▶ recursos de la Linguist List
- ▶ WEBIR/IE recursos de IR, publicaciones, conferencias, contactos, listas de noticias...

instituciones

- ▶ **ELDA - ELRA** Evaluations and Language resources Distribution Agency, tiene un completísimo **catálogo de recursos lingüísticos** para lenguas europeas, libres y pagos.
- ▶ **HLT central** Human Language Technology, repositorio europeo de grupos y entidades relacionados con las tecnologías del lenguaje, tiene un completo calendario de eventos y enlaces interesantes
- ▶ **ELSNET** European Network of Excellence in Human Language Technologies, con calendario de eventos (es uno de los principales sponsors del área), grupos relacionados, asociaciones, bolsa de trabajo y de becas, etc.
- ▶ **ACL** The Association for Computational Linguistics tiene enlaces a las principales conferencias mundiales sobre lenguaje

empresas

- ▶ **MITRE** tiene muchos proyectos de investigación en lenguaje natural, muchos con recursos libres
- ▶ **Xerox**
- ▶ **AT&T**