

# Introducción al Procesamiento de Lenguaje Natural

Dr. Gabriel Infante-Lopez  
Sección Informática - FaMAF  
Universidad Nacional de Córdoba  
[gabriel@famaf.unc.edu.ar](mailto:gabriel@famaf.unc.edu.ar)  
<http://www.cs.famaf.unc.edu.ar/~gabriel>

# Material de Lectura

- Manning, C., Schütze, H. *Foundations of Statistical Natural Language Processing*. The MIT press. 1999. **Obligatorio**
- Jurafsky, D., Martin, J. *Speech and Language Processing*. Prentice-Hall. 2000. Referencia
- Wall, L. et al. *Programming Perl*. 3rd Ed. O'Reilly. Referencia
- Proceedings de conferencias: ACL, EACL, COLING

# Requirimientos del Curso pesos para la nota final

- Proyectos Individuales (4) 40%
- Examen parcial mitad cuatrimestre 30%
- Examen parcial final cuatrimestre 30%
- Los exámenes deben ser aprobados para que promediar los resultados.

# Proyectos y Organización

## ● Proyectos

- Entropía, y modelos de lenguaje

- Clases de Palabra

- Clasificación

- Sintaxis

## ● Organización

- algunos ejercicios de papel y lapiz, mucha programación

- deadlines estrictos. 10 días de tardanza máximo. Ultimo proyecto solo 2 días.

- Plagio NO permitido



# Descripción del Curso

- Intro., Probabilidad, y teoría de la información
  - Conceptos básicos: definiciones, formulas, ejemplo, etc.
- Modelos de Lenguaje
  - n-grams, estimación de parámetros
  - smoothing (EM algorithm)
- Un poco de Lingüística
  - fonología, morfología, sintaxis, semántica, discurso
- Palabras y Léxico
  - clases de palabra, información mutua

# Descripción del Curso

- Hidden Markov Models
  - background, algoritmos, estimación param.
- Tagging: Metodos, Algoritmos, Evaluación
  - Tagsets, morfología, lematización
  - HMM tagging, Transformation y feature based
- Gramaticas, y Analisis Sintactico: Datos, Algs.
  - Gramaticas y automatas, parsing determ.
- Aplicaciones.

# PLN: Principales problemas

- varias palabras, varios fenomenos  $\Rightarrow$  alta complejidad. Ejemplo: Lexicón Finlandes  $\sim 2 \times 10^7$
- irregularidades (excepción, excepción de la excepción):
  - cantar  $\rightarrow$  canté, hablar  $\rightarrow$  hablé, estar  $\rightarrow$  estuve.

# PLN: Principales problemas

## ● Ambigüedad:

● *bajo*: Adjetivo, nombre, preposición

● *Prohibido girar a la izquierda de 12 a 14. Excepto taxis y remises. ¿Los taxis y remises no pueden girar nunca?*

● *Los chorros de aguas cordobesas. ¿Cuántas interpretaciones posibles?*

● *Gracias por no comer, beber y escuchar música sin auriculares. ¿Podemos comer con auriculares?*



# ¿Reglas categoricas o Estadística?

- Preferencias:
  - casos claros: pistas en el contexto: *bajo la mesa* (bajo es una preposición)
  - casos menos claros. *El viento le voló el sombrero, el hombre lo quiso agarrar. A quien? ¿al sombrero o al viento?*
  - semanticos: *comer pizza con cuchillo, comer pizzas con aceitunas.*

# Soluciones

- Si sabemos, mejor no adivinar
  - Morfología, inflexiones, excepciones
  - lexicons, lista de palabras
  - nombres inambiguos
  - frases fijas, colocaciones: ej. Aguas Cordobesa.
  - Reglas sintacticas?  $S \rightarrow Np Vp$
- Usar estadísticas solo para preferencias

# NLP Estadístico

- Supongamos que a cada oración  $W$ , se le asigna una probabilidad  $P(W|X)$  en un contexto  $X$ .
- Para cada contexto  $X$ , ordenar todas las oraciones de acuerdo a  $P(W|X)$ .

