

¿Qué es lingüística computacional y cómo se relaciona con el procesamiento de lenguaje natural?

La **lingüística computacional** es un campo interdisciplinar que combina la informática y la lingüística para estudiar cómo las máquinas pueden procesar, entender y generar lenguaje humano. En el libro se introduce a través de la evolución de sistemas como **ELIZA**, que imitaba a un terapeuta Rogeriano usando patrones de texto simples. Esto marca el inicio del **procesamiento de lenguaje natural (PLN)**, el cual va más allá de simples coincidencias de texto para abarcar tareas como segmentación, análisis sintáctico y generación de lenguaje, usando algoritmos y modelos estadísticos.

¿Cuáles son y en qué consisten cada uno de los componentes para el uso eficiente de textos?

El libro describe varios **componentes fundamentales del preprocesamiento del lenguaje**:

1. **Tokenización:** Segmentación de texto en palabras o subunidades, considerando complejidades como emoticones o contracciones.
2. **Normalización de texto:** Incluye:
 - **Lematización:** Reducir palabras a su forma base (ej. *sings, sang* → *sing*).
 - **Stemming:** Versión más simple que corta sufijos sin analizar morfología.
 - **Case folding:** Convertir a minúsculas.
3. **Segmentación de oraciones:** Identificar límites de oración basados en puntuación, considerando abreviaciones.
4. **Cálculo de distancias de edición:** Medida para evaluar similitud entre cadenas (útil en corrección ortográfica, reconocimiento de voz, etc.).
5. **Uso de expresiones regulares:** Herramienta poderosa para encontrar patrones en texto (fechas, precios, frases específicas).
6. **Tokenización por subpalabras (BPE):** Divide palabras en fragmentos frecuentes para manejar palabras raras o nuevas en grandes modelos de lenguaje.

7. **Corpus y variación lingüística:** Consideraciones sobre diversidad de lenguas, géneros, dialectos, y fenómenos como el *code-switching*.
-

¿Qué problemas acarrear los ‘ayudantes naturales’?

El texto menciona a ELIZA como un ejemplo temprano de ‘ayudante natural’. Aunque era percibido como inteligente, **no entendía realmente el lenguaje**, sino que se basaba en reglas fijas de sustitución y coincidencia de patrones. El problema con estos sistemas es que:

- Son superficiales y no comprenden el significado real.
 - Generan falsas expectativas en los usuarios.
 - No generalizan ni manejan contextos o ambigüedades reales del lenguaje humano.
-

De las aplicaciones presentadas, ¿cuál es la que más te interesa y cómo la aplicarías en tu contexto laboral/académico?

Una de las aplicaciones más destacadas es el uso de **tokenización por subpalabras y lematización** para crear representaciones normalizadas del lenguaje, fundamentales para modelos estadísticos y de aprendizaje automático. En un contexto académico o laboral enfocado en análisis textual (como minería de opiniones o clasificación de documentos), estas técnicas permiten:

- Preprocesar grandes volúmenes de texto de forma eficiente.
- Mejorar la calidad de modelos de clasificación o análisis de sentimientos.
- Reducir la dimensionalidad del vocabulario manejado por los modelos.

En particular, si trabajas con análisis de contenido o datos en redes sociales, usar BPE y lematización puede ayudarte a manejar con precisión textos informales y multilingües.