



# Profundización del Conocimiento

## Tokenización en el PLN

La tokenización es el proceso de segmentar un texto en unidades más pequeñas denominadas tokens. Un token puede ser una palabra, una frase, un carácter o incluso una subpalabra, dependiendo del enfoque utilizado. Este proceso es una de las primeras etapas en el Procesamiento del Lenguaje Natural (PLN), ya que permite estructurar el texto de manera que pueda ser analizado y manipulado por algoritmos y modelos computacionales.

El propósito de la tokenización es dividir el texto en fragmentos manejables que conserven significado lingüístico. Por ejemplo, en una oración como:

*"El aprendizaje profundo ha revolucionado el PLN."*

La tokenización a nivel de palabras resultaría en:

["El", "aprendizaje", "profundo", "ha", "revolucionado", "el", "PLN", "."]

Sin embargo, dependiendo del idioma y del contexto, la tokenización puede no ser trivial, especialmente en lenguas que no utilizan espacios para separar palabras, como el chino, el japonés y el tailandés.

### Tipos de Tokenización

Existen varios enfoques para la tokenización, cada uno con sus ventajas y desafíos. Entre los más comunes se encuentran:

#### *Tokenización por Palabras*

Este método segmenta el texto en palabras individuales. Se utiliza en lenguajes que tienen delimitadores claros, como el español e inglés.

Ejemplo: *"La inteligencia artificial está avanzando rápidamente."*

Tokens: ["La", "inteligencia", "artificial", "está", "avanzando", "rápidamente", "."]

#### *Tokenización por Caracteres*

Cada carácter se trata como un token. Se usa en aplicaciones específicas, como modelos de generación de texto o análisis de idiomas como el chino y japonés.

Ejemplo: *"Hola"*

Tokens: ["H", "o", "l", "a"]

#### *Tokenización por Subpalabras*

Este método segmenta palabras en unidades más pequeñas basadas en su morfología. Es útil en PLN moderno para manejar palabras raras o desconocidas. Métodos como Byte Pair Encoding (BPE) y WordPiece son ampliamente utilizados en modelos como GPT y BERT.



# Profundización del Conocimiento

Ejemplo con BPE:

"inalterablemente" → ["in", "alter", "able", "mente"]

## Tokenización por Frases

Segmenta el texto en frases o en partes que mantienen sentido completo. Es útil en tareas como la traducción automática.

Ejemplo:

"El sol brilla. Hace calor."

Tokens: ["El sol brilla.", "Hace calor."]

## Importancia de la Tokenización en PLN

La tokenización es crucial en diversas aplicaciones de PLN, como:

- **Análisis de Sentimiento:** Divide el texto para identificar palabras clave con connotaciones positivas o negativas.
- **Traducción Automática:** Facilita la segmentación del texto en unidades que pueden ser traducidas con mayor precisión.
- **Reconocimiento de Entidades:** Ayuda a detectar nombres de personas, lugares y organizaciones en textos.
- **Modelado del Lenguaje:** Modelos como GPT-4 o BERT dependen de una tokenización eficiente para aprender patrones lingüísticos.

## Desafíos en la Tokenización

A pesar de ser un paso fundamental en PLN, la tokenización presenta desafíos, tales como:

- **Ambigüedad léxica:** En frases como "*El banco está cerrado*", la palabra *banco* puede referirse a una entidad financiera o a un asiento.
- **Contracciones:** En inglés, términos como "*can't*" requieren decisiones sobre si dividirlos en ["*can*", "*not*"] o tratarlos como un único token.
- **Lenguas sin espacios:** Idiomas como el chino no usan espacios entre palabras, lo que dificulta la segmentación.
- **Ortografía y errores tipográficos:** Palabras mal escritas pueden dificultar la tokenización correcta.



# Profundización del Conocimiento

## Herramientas de Tokenización

Existen varias herramientas y bibliotecas diseñadas para tokenizar texto de manera eficiente, entre ellas:

- NLTK (Natural Language Toolkit): Popular en Python para tokenización de textos en múltiples idiomas.
- spaCy: Biblioteca avanzada que ofrece tokenización rápida y precisa.
- Tokenizers de Hugging Face: Implementaciones optimizadas para modelos de aprendizaje profundo.
- Mecab y Jieba: Herramientas diseñadas para tokenizar textos en japonés y chino, respectivamente.

## La pragmática en el PLN

La pragmática es una rama de la lingüística que estudia el uso del lenguaje en contexto y cómo los hablantes interpretan los significados más allá de la estructura gramatical y semántica de las palabras. Es decir, se centra en la relación entre el lenguaje y sus usuarios en situaciones comunicativas específicas. La pragmática analiza cómo el contexto influye en la interpretación del significado, considerando elementos como el tono, la intención del hablante, la ironía, el sarcasmo y otros factores que no pueden ser interpretados solo mediante la estructura lingüística.

Por ejemplo, en la frase:

*"Hace frío aquí."*

El significado pragmático puede variar según el contexto:

- Puede ser una simple observación sobre la temperatura.
- Puede ser una sugerencia indirecta para cerrar una ventana.
- Puede expresar una queja o una petición implícita.

Este análisis ilustra que la pragmática va más allá del significado literal y explora cómo el lenguaje se usa en interacciones reales.

## Características Principales de la Pragmática

La pragmática se distingue por los siguientes aspectos fundamentales:

1. Contexto: La pragmática depende del entorno en el que se usa el lenguaje, lo que incluye el lugar, el momento, la relación entre los hablantes y sus intenciones comunicativas.



# Profundización del Conocimiento

2. Intención del hablante: No siempre lo que se dice explícitamente refleja el verdadero significado. El lenguaje suele ser usado para transmitir significados implícitos.
3. Inferencias y supuestos compartidos: Se basa en la capacidad del oyente para hacer inferencias a partir del contexto y del conocimiento previo compartido con el hablante.
4. Polisemia e ironía: La pragmática ayuda a interpretar expresiones con múltiples significados o aquellas con intención humorística o sarcástica.

## Principales Teorías de la Pragmática

A lo largo del tiempo, varios lingüistas han desarrollado teorías para explicar cómo funciona la pragmática en la comunicación humana. Algunas de las más influyentes son:

### *Principio de Cooperación de Grice*

Paul Grice (1975) propuso que la comunicación efectiva se basa en un principio de cooperación, donde los interlocutores asumen que ambos buscan comprenderse mutuamente. Para ello, formuló las cuatro máximas de Grice:

- Máxima de cantidad: Se debe proporcionar la cantidad justa de información, ni más ni menos de lo necesario.
- Máxima de calidad: Se espera que la información sea veraz y confiable.
- Máxima de relación: Lo que se dice debe ser relevante para la conversación.
- Máxima de manera: Se debe hablar de forma clara y ordenada.

Por ejemplo, si alguien pregunta:

*"¿Qué hora es?"*

Responder *"Las manecillas del reloj marcan una posición determinada."* violaría la máxima de manera, ya que es una respuesta innecesariamente complicada.

### *Teoría de los Actos de Habla de Austin y Searle*

J.L. Austin (1962) y John Searle (1969) introdujeron la teoría de los actos de habla, que clasifica las expresiones lingüísticas según su función comunicativa:

- Actos locutivos: Lo que se dice literalmente.
- Actos ilocutivos: La intención con la que se dice algo.
- Actos perlocutivos: El efecto que la expresión tiene en el oyente.



# Profundización del Conocimiento

Ejemplo:

Si alguien dice *"Te prometo que llegaré temprano"*, no solo está afirmando algo, sino que está realizando el acto de hacer una promesa.

*Teoría de la Relevancia de Sperber y Wilson*

Dan Sperber y Deirdre Wilson (1986) propusieron que la comunicación sigue un principio de relevancia, según el cual los hablantes estructuran su mensaje para maximizar su comprensión con el menor esfuerzo cognitivo posible. Esto explica por qué ciertas expresiones implícitas o indirectas pueden ser más efectivas que una explicación detallada.

Ejemplo:

Si alguien pregunta: *"¿Vamos al cine?"* y el otro responde *"Tengo un examen mañana"*, la inferencia pragmática es que no puede ir al cine, aunque no lo haya dicho explícitamente.

## Aplicaciones de la Pragmática

La pragmática tiene múltiples aplicaciones en diversas disciplinas, desde la lingüística hasta la inteligencia artificial.

*Pragmática en la Lingüística Aplicada*

En la enseñanza de idiomas, comprender la pragmática es crucial para que los hablantes no solo aprendan gramática y vocabulario, sino también las normas culturales de comunicación. Por ejemplo, en inglés, rechazar una invitación de manera directa puede ser visto como descortés, mientras que en español es más común.

Ejemplo:

- En inglés: *"Sorry, I can't go."*
- En español: *"Me encantaría, pero tengo otro compromiso."* (una forma más atenuada)

*Pragmática en el Procesamiento del Lenguaje Natural (PLN)*

Los sistemas de inteligencia artificial y asistentes virtuales (como Siri, Alexa o ChatGPT) dependen de modelos pragmáticos para interpretar correctamente el significado de las consultas de los usuarios.

Ejemplo:

Si un usuario le dice a un asistente virtual:

*"Tengo frío."*

Un sistema pragmáticamente avanzado puede sugerir subir la calefacción, en lugar de responder con una definición sobre la temperatura.



# Profundización del Conocimiento

## *Pragmática en la Comunicación Intercultural*

Las diferencias culturales influyen en cómo se usa el lenguaje. Un gesto o expresión que en una cultura es normal, en otra puede ser ofensivo. Por ejemplo, en Japón, el silencio en una conversación puede significar respeto y reflexión, mientras que en países occidentales puede interpretarse como incomodidad o falta de interés.

## Desafíos en la Pragmática

A pesar de su importancia, la pragmática enfrenta varios desafíos en su estudio y aplicación:

1. Ambigüedad e ironía: La ironía y el sarcasmo pueden ser difíciles de detectar, incluso para hablantes nativos.
2. Falta de marcadores formales: Mientras que la gramática y la semántica tienen estructuras definidas, la pragmática depende más de la intuición y del contexto.
3. Procesamiento por máquinas: Aunque los modelos de IA han mejorado en la interpretación pragmática, aún hay desafíos en la comprensión de humor, dobles sentidos y referencias culturales.