# Reflection Journal: Lab 02 - Basic NLP Preprocessing Techniques

**Marvin Azuougu**

**Department of Science, Technology, Engineering & Math**

**Houston Community College**

**ITAI-2373-Natural Language Processing Summer 2025**

**Professor: Patricia McManus**

**10th of June 2025.**

This lab on basic NLP preprocessing was an illuminating experience that deepened my understanding of what truly makes natural language processing "intelligent." Before diving in, preprocessing felt like a mere preliminary step—a necessary chore before the real work began. However, this module emphatically drove home its pivotal role in shaping the effectiveness, efficiency, and even the interpretability of any downstream NLP task.

One of my most significant insights was grasping the sheer "messiness" of human language in its raw form. The initial conceptual questions, prompting me to consider the inherent challenges faced by systems like search engines and chatbots, immediately resonated. Ambiguity, nuance (including sarcasm and irony), and the dynamic, informal nature of daily communication, complete with slang, emojis, and misspellings, present complexities far beyond what I initially appreciated. Preprocessing isn't simply about removing noise; it's about systematically managing this linguistic chaos to present a standardized, digestible format for algorithms. The stark contrast between raw social media text and its preprocessed counterpart, often stripping away valuable context for generalization, vividly illustrated this point.

Working through the various tokenization approaches, particularly the comparison between NLTK and spaCy, was a turning point. While NLTK's word tokenize provided a foundational understanding of breaking text into units, spaCy's integrated pipeline felt like a monumental leap. The fact that spaCy immediately provides Part-of-Speech (POS) tags, lemmas, and even dependency parses right after tokenization underscores its design philosophy for robust, production-ready applications. This was a true "lightbulb moment" for me; I realized that effective preprocessing isn't just a sequence of isolated steps, but a holistic process where components seamlessly build upon one another. My initial challenge of grasping how these layers contribute quickly dissipated as I observed spaCy's unified output, leading me to ponder if NLTK's more modular approach, while flexible, might introduce unnecessary complexity for many standard tasks.

The segment on stop word removal and the detailed comparison between NLTK's and spaCy's lists was surprisingly thought-provoking. I had always assumed stop words were universally agreed upon, but the subtle differences in their curated lists highlighted a crucial trade-off: noise reduction versus information preservation. For instance, removing intensifiers like "very" or "really" might seem like an obvious noise reduction step, yet it can severely undermine sentiment analysis by stripping away crucial emotional intensity. This directly fed into a core question that arose for me: How does one judiciously decide which words are truly "stop" words for a specific domain or task, especially when context is paramount? The observation that spaCy's more conservative stop word list, sometimes

retaining these subtle but significant words, could enhance accuracy in nuanced applications was a key takeaway.

The detailed comparison of stemming versus lemmatization further amplified this fundamental trade-off. Stemming's aggressive, rule-based approach, often yielding non-dictionary words (like "bett" from "better") and failing to normalize irregular forms, felt almost crude compared to lemmatization's context-aware, linguistically sound method that returns valid dictionary forms (like "good" from "better"). This distinction immediately connected to real-world applications. For a search engine, where speed and high recall are primary drivers, stemming might offer a sufficient, quick solution. However, for a sentiment analysis system or a real-time chatbot, where meaning preservation and accuracy are non-negotiable, lemmatization emerges as the unequivocal choice. The loss of subtle semantic links from aggressive stemming, such as "better" not becoming "good," could be disastrous for models striving to grasp genuine intent and nuance.

My exploration also raised crucial questions about the optimal balance for different text types. It became evident that social media text was the most profoundly affected by preprocessing due to its inherently informal and unconventional linguistic features. Conversely, academic text was notably the least affected, largely owing to its formal and highly structured nature. This variability underscores that there's no singular "best" preprocessing pipeline, emphasizing the critical need for adaptive strategies tailored to the specific characteristics of the data.

Looking ahead, the insights gleaned from this lab will be invaluable. If I were to build an NLP system to analyze customer reviews, I would unequivocally opt for a Standard Processing pipeline augmented with lemmatization, prioritizing accuracy and meaning preservation. My approach would involve meticulously handling stop words, ensuring intensifiers aren't discarded, and, critically, incorporating explicit processing for emojis and impactful punctuation (like multiple exclamation marks or ellipses), as they are direct indicators of sentiment in informal review contexts. I would deliberately avoid aggressive stemming, recognizing that the nuances of customer feedback are far too significant to sacrifice. This experience has reinforced that the central trade-offs in any NLP project—accuracy vs. speed, information preservation vs. noise reduction, and generalizability vs. specificity—are not theoretical dilemmas but practical decisions that directly shape a system's efficacy. This lab wasn't just about mastering techniques; it was about truly understanding the intricate art and science of preparing text for intelligent analysis.