

What are Language Models in NLP?

Language models are a fundamental component of natural language processing (NLP) and computational linguistics. They are designed to understand, generate, and predict human language. These models analyze the structure and use of language to perform tasks such as machine translation, text generation, and sentiment analysis.

What is a Language Model in Natural Language Processing?

A language model in natural language processing (NLP) is a statistical or machine learning model that is used to predict the next word in a sequence given the previous words. Language models play a crucial role in various NLP tasks such as machine translation, speech recognition, text generation, and sentiment analysis. They analyze and understand the structure and use of human language, enabling machines to process and generate text that is contextually appropriate and coherent.

Language models can be broadly categorized into two types:

- Pure Statistical Methods
- Neural Models

Purpose and Functionality

The primary purpose of a language model is to capture the statistical properties of natural language. By learning the probability distribution of word sequences, a language model can predict the likelihood of a given word following a sequence of words. This predictive capability is fundamental for tasks that require understanding the context and meaning of text.

For instance, in text generation, a language model can generate plausible and contextually relevant text by predicting the next word in a sequence iteratively. In machine translation, language models help in translating text from one language to another by understanding and generating grammatically correct sentences in the target language.

Pure Statistical Methods

Pure statistical methods form the basis of traditional language models. These methods rely on the statistical properties of language to predict the next word in a sentence, given the previous words. They include n-grams, exponential models, and skip-gram models.

1. N-grams

An n-gram is a sequence of n items from a sample of text or speech, such as phonemes, syllables, letters, words, or base pairs. N-gram models use the frequency of these sequences in a training corpus to predict the likelihood of word sequences. For example, a bigram (2-gram) model predicts the next word based on the previous word, while a trigram (3-gram) model uses the two preceding words.

N-gram models are simple, easy to implement, and computationally efficient, making them suitable for applications with limited computational resources. However, they have significant limitations. They struggle with capturing long-range dependencies due to their limited context window. As n increases, the number of possible n-grams grows exponentially, leading to sparsity issues where many sequences are never observed in the training data. This sparsity makes it difficult to accurately estimate the probabilities of less common sequences.

2. Exponential Models

Exponential models, such as the Maximum Entropy model, are more flexible and powerful than n-gram models. They predict the probability of a word based on a wide range of features, including not only the previous words but also other contextual information. These models assign weights to different features and combine them using an exponential function to estimate probabilities.

Maximum Entropy Models

Maximum Entropy (MaxEnt) models, also known as logistic regression in the context of classification, are used to estimate the

probabilities of different outcomes based on a set of features. In the context of language modeling, MaxEnt models use features such as the presence of certain words, part-of-speech tags, and syntactic patterns to predict the next word. The model parameters are learned by maximizing the likelihood of the observed data under the model.

MaxEnt models are more flexible than n-gram models because they can incorporate a wider range of features. However, they are also more complex and computationally intensive to train. Like n-gram models, MaxEnt models still struggle with long-range dependencies because they rely on fixed-length context windows.

3. Skip-gram Models

Skip-gram models are a type of statistical method used primarily in word embedding techniques. They predict the context words (surrounding words) given a target word within a certain window size. Skip-gram models, particularly those used in Word2Vec, are effective for capturing the semantic relationships between words by optimizing the likelihood of context words appearing around a target word.

Word2Vec and Skip-gram

Word2Vec, developed by Google, includes two main architectures: skip-gram and continuous bag-of-words (CBOW). The skip-gram model predicts the context words given a target word, while the CBOW model predicts the target word given the context words. Both models are trained using neural networks, but they are conceptually simple and computationally efficient.

Neural Models

Neural models have revolutionized the field of NLP by leveraging deep learning techniques to create more sophisticated and accurate language models. These models include Recurrent Neural Networks (RNNs), Transformer-based models, and large language models.

1. Recurrent Neural Networks

Recurrent Neural Networks (RNNs) are a type of neural network designed for sequential data, making them well-suited for language modeling. RNNs maintain a hidden state that captures information about previous inputs, allowing them to consider the context of words in a sequence.

LSTMs and GRUs are advanced RNN variants that address the vanishing gradient problem, enabling the capture of long-range dependencies in text. LSTMs use a gating mechanism to control the flow of information, while GRUs simplify the gating mechanism, making them faster to train.

2. Transformer-based Models

The Transformer model, introduced by Vaswani et al. in 2017, has revolutionized NLP. Unlike RNNs, which process data sequentially, the Transformer model processes the entire input simultaneously, making it more efficient for parallel computation.

The key components of the Transformer architecture are:

Self-Attention Mechanism: This mechanism allows the model to weigh the importance of different words in a sequence, capturing dependencies regardless of their distance in the text. Each word's representation is updated based on its relationship with all other words in the sequence.

Encoder-Decoder Structure: The Transformer consists of an encoder and a decoder. The encoder processes the input sequence and generates a set of hidden representations. The decoder takes these representations and generates the output sequence.

Positional Encoding: Since Transformers do not process the input sequentially, they use positional encoding to retain information about the order of words in a sequence. This encoding adds positional information to the input embeddings, allowing the model to consider the order of words.

Some of the transformers based models are - BERT, GPT-3, T5 and more.

3. Large Language Models (LLMs)

Large language models have pushed the boundaries of what is possible in NLP. These models are characterized by their vast size, often comprising billions of parameters, and their ability to perform a wide range of tasks with minimal fine-tuning.

Training large language models involves feeding them vast amounts of text data and optimizing their parameters using powerful computational resources. The training process typically includes multiple stages, such as unsupervised pre-training on large corpora followed by supervised fine-tuning on specific tasks.

While large language models offer remarkable performance, they also pose significant challenges. Training these models requires substantial computational resources and energy, raising concerns about their environmental impact. Additionally, the models' size and complexity can make them difficult to interpret and control, leading to potential ethical and bias issues.

Popular Language Models in NLP

Several language models have gained prominence due to their innovative architecture and impressive performance on NLP tasks.

Here are some of the most notable models:

BERT (Bidirectional Encoder Representations from Transformers)

BERT, developed by Google, is a Transformer-based model that uses bidirectional context to understand the meaning of words in a sentence. It has improved the relevance of search results and achieved state-of-the-art performance in many NLP benchmarks.

GPT-3 (Generative Pre-trained Transformer 3)

GPT-3, developed by OpenAI, is a large language model known for its ability to generate coherent and contextually appropriate text based on a given prompt. With 175 billion parameters, it is one of the largest and most powerful language models to date.

T5 (Text-to-Text Transfer Transformer)

T5, developed by Google, treats all NLP tasks as a text-to-text problem, enabling it to handle a wide range of tasks with a single model. It has demonstrated versatility and effectiveness across various NLP tasks.

Word2Vec

Word2Vec, developed by Google, includes the skip-gram and continuous bag-of-words (CBOW) models. These models create word embeddings that capture semantic similarities between words, improving the performance of downstream NLP tasks.

ELMo (Embeddings from Language Models)

ELMo generates context-sensitive word embeddings by considering the entire sentence. It uses bidirectional LSTMs and has improved performance on various NLP tasks by providing more nuanced word representations.

Transformer-XL

Transformer-XL is an extension of the Transformer model that addresses the fixed-length context limitation by introducing a segment-level recurrence mechanism. This allows the model to capture longer-range dependencies more effectively.

XLNet

XLNet, developed by Google, is an autoregressive Transformer model that uses permutation-based training to capture bidirectional context. It has achieved state-of-the-art results on several NLP benchmarks.

RoBERTa (Robustly Optimized BERT Approach)

RoBERTa, developed by Facebook AI, is a variant of BERT that uses more extensive training data and optimizations to achieve better performance. It has set new benchmarks in several NLP tasks.

ALBERT (A Lite BERT)

ALBERT, developed by Google, is a lightweight version of BERT that reduces the model size while maintaining performance. It achieves this by sharing parameters across layers and factorizing the embedding parameters.

Turing-NLG

Turing-NLG, developed by Microsoft, is a large language model known for its ability to generate high-quality text. It has been used in various applications, including chatbots and virtual assistants.

N-Gram Language Modelling with NLTK

Language modeling is the way of determining the probability of any sequence of words. Language modeling is used in various applications such as Speech Recognition, Spam filtering, etc. Language modeling is the key aim behind implementing many state-of-the-art Natural Language Processing models.

Methods of Language Modelling

Two methods of Language Modeling:

Statistical Language Modelling: Statistical Language Modeling, or Language Modeling, is the development of probabilistic models that can predict the next word in the sequence given the words that precede. Examples such as N-gram language modeling.

Neural Language Modeling: Neural network methods are achieving better results than classical methods both on standalone language models and when models are incorporated into larger models on challenging tasks like speech recognition and machine translation. A way of performing a neural language model is through word embeddings.

N-gram

N-gram can be defined as the contiguous sequence of n items from a given sample of text or speech. The items can be letters, words, or base pairs according to the application. The N-grams typically are collected from a text or speech corpus (A long text dataset).

For instance, N-grams can be unigrams like (“This”, “article”, “is”, “on”, “NLP”) or bigrams (“This article”, “article is”, “is on”, “on NLP”).

N-gram Language Model

An N-gram language model predicts the probability of a given N-gram within any sequence of words in a language. A well-crafted N-gram model can effectively predict the next word in a sentence, which is essentially determining the value of $p(w|h)$, where h is the history or context and w is the word to predict.

Let's explore how to predict the next word in a sentence. We need to calculate $p(w|h)$, where w is the candidate for the next word. Consider the sentence ‘This article is on...’. If we want to calculate the probability of the next word being “NLP”, the probability can be expressed as:

$$p(\text{“NLP”}|\text{“This”, “article”, “is”, “on”})p(\text{“NLP”}|\text{“This”, “article”, “is”, “on”})$$

To generalize, the conditional probability of the fifth word given the first four can be written as:

$$p(w_5|w_1, w_2, w_3, w_4) \text{ or } p(W) = p(w_n|w_1, w_2, \dots, w_{n-1}) \text{ or } p(W) = p(w_n|w_1, w_2, \dots, w_{n-1})$$

This is calculated using the chain rule of probability:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \text{ and } P(A \cap B) = P(A|B)P(B) \Rightarrow P(A|B) = \frac{P(A \cap B)}{P(B)} \text{ and } P(A \cap B) = P(A|B)P(B)$$

Now generalize this to sequence probability:

$$P(X_1, X_2, \dots, X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2) \dots P(X_n|X_1, X_2, \dots, X_{n-1})$$
$$P(X_1, X_2, \dots, X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2) \dots P(X_n|X_1, X_2, \dots, X_{n-1})$$

This yields:

$$P(w_1, w_2, w_3, \dots, w_n) = \prod_i P(w_i|w_1, w_2, \dots, w_{i-1}) \Rightarrow P(w_1, w_2, w_3, \dots, w_n) = \prod_i P(w_i|w_1, w_2, \dots, w_{i-1})$$

By applying Markov assumptions, which propose that the future state depends only on the current state and not on the sequence of events that preceded it, we simplify the formula:

$$P(w_i|w_1, w_2, \dots, w_{i-1}) \approx P(w_i|w_{i-k}, \dots, w_{i-1}) \Rightarrow P(w_i|w_1, w_2, \dots, w_{i-1}) \approx P(w_i|w_{i-k}, \dots, w_{i-1})$$

For a unigram model ($k=0$), this simplifies further to:

$$P(w_1, w_2, \dots, w_n) \approx \prod_i P(w_i) \Rightarrow P(w_1, w_2, \dots, w_n) \approx \prod_i P(w_i)$$

And for a bigram model ($k=1$):

$$P(w_i|w_1, w_2, \dots, w_{i-1}) \approx P(w_i|w_{i-1}) \Rightarrow P(w_i|w_1, w_2, \dots, w_{i-1}) \approx P(w_i|w_{i-1})$$

Metrics for Language Modelling

- **Entropy:** Entropy, as a measure of the amount of information conveyed by Claude Shannon. Below is the formula for representing entropy

$$H(p) = -\sum_x p(x) \cdot \log(p(x)) \quad H(p) = -\sum_x p(x) \cdot \log(p(x))$$

$H(p)$ is always greater than or equal to 0.

- **Cross-Entropy:** It measures the ability of the trained model to represent test data ($w_{1i-1} \dots w_{1i-1}$).

$$H(p) = -\sum_{i=1}^n \log_2(p(w_i | w_{1i-1})) \quad H(p) = -\sum_{i=1}^n \log_2(p(w_i | w_{1i-1}))$$

The cross-entropy is always greater than or equal to Entropy i.e the model uncertainty can be no less than the true uncertainty.

- **Perplexity:** Perplexity is a measure of how good a probability distribution predicts a sample. It can be understood as a measure of uncertainty. The perplexity can be calculated by cross-entropy to the exponent of 2.

$$2^{\text{Cross-Entropy}}$$

Following is the formula for the calculation of Probability of the test set assigned by the language model, normalized by the number of words:

$$PP(W) = \prod_{i=1}^n P(w_i | w_{i-1}) \quad PP(W) = \prod_{i=1}^n P(w_i | w_{i-1})$$

For Example:

- Let's take an example of the sentence: '**Natural Language Processing**'. For predicting the first word, let's say the word has the following probabilities:

word	P(word <start>)
The	0.4
Processing	0.3
Natural	0.12
Language	0.18

- Now, we know the probability of getting the first word as natural. But, what's the probability of getting the next word after getting the word '*Language*' after the word '*Natural*'.

word	P(word 'Natural')
The	0.05
Processing	0.3
Natural	0.15
Language	0.5

- After getting the probability of generating words 'Natural Language', what's the probability of getting '*Processing*'.

word	P(word 'Language')
The	0.1
Processing	0.7
Natural	0.1
Language	0.1

- Now, the perplexity can be calculated as:

$$PP(W) = \prod_{i=1}^N P(w_i | w_{i-1}) = 0.12 * 0.5 * 0.73 \approx 2.876$$

$$PP(W) = \prod_{i=1}^n P(w_i | w_{i-1}) = 0.12 * 0.5 * 0.71 \approx 2.876$$

- From that we can also calculate entropy:

$$Entropy = \log_2(2.876) = 1.524$$

$$Entropy = \log_2(2.876) = 1.524$$

Shortcomings:

- To get a better context of the text, we need higher values of n, but this will also increase computational overhead.
- The increasing value of n in n-gram can also lead to sparsity.

What is smoothing in NLP and why do we need it

What is smoothing in the context of natural language processing, define smoothing in NLP, what is the purpose of smoothing in nlp, is smoothing an important task in language model

Smoothing in NLP

Smoothing is the process of flattening a probability distribution implied by a language model so that all reasonable word sequences can occur with some probability. This often involves broadening the distribution by redistributing weight from high probability regions to zero probability regions.

Smoothing not only prevents zero probabilities, attempts to improve the accuracy of the model as a whole.

Why do we need smoothing?

In a language model, we use parameter estimation (MLE) on training data. We can't actually evaluate our MLE models on unseen test data because both are likely to contain words/n-grams that these models assign zero probability to. Relative frequency estimation assigns all probability mass to events in the training corpus. But we need to reserve some probability mass to events that don't occur (unseen events) in the training data.

Example:

Training data: *The cow is an animal.*

Test data: *The dog is an animal.*

If we use **unigram model** to train;

$P(\text{the}) = \text{count}(\text{the}) / (\text{Total number of words in training set}) = 1/5.$

Likewise, $P(\text{cow}) = P(\text{is}) = P(\text{an}) = P(\text{animal}) = 1/5$

To evaluate (test) the **unigram model**;

$P(\text{the cow is an animal}) = P(\text{the}) * P(\text{cow}) * P(\text{is}) * P(\text{an}) * P(\text{animal}) = 0.00032$

While we use unigram model on the test data, it becomes zero because $P(\text{dog}) = 0$. The term 'dog' never occurred in the training data. Hence, we use smoothing.

Advanced Smoothing Techniques in Language Models

Language models predict the probability of a sequence of words and generate coherent text. These models are used in various applications, including chatbots, translators, and more. However, one of the challenges in building language models is handling the issue of zero probabilities for unseen events in the training data. Smoothing techniques are employed to address this problem, ensuring that predictions remain accurate even when encountering previously unseen words.

Understanding Language Models

Large Language Models (LLMs) are designed to understand and generate human-like text. They work by calculating the probability of a sequence of words and predicting the next word in a sentence. These models are typically trained on vast corpora of text, enabling them to grasp the patterns and structures of language.

LLMs, like those utilized in AI applications, often use deep learning architectures such as Recurrent Neural Networks (RNNs) and Transformers. These models can capture long-range dependencies and complex linguistic patterns, making them highly effective for a wide range of NLP tasks.

Smoothing techniques are crucial in large language models to tackle the problem of zero probabilities for unseen words. These techniques adjust the probability distribution, ensuring that even unseen events have a non-zero probability. By doing so, they significantly enhance the performance of LLMs.

1. Witten-Bell Smoothing

Witten-Bell smoothing is a technique used in statistical language modeling to estimate the probability of unseen events, particularly in the context of n-grams in natural language processing (NLP). This method was

introduced by Ian H. Witten and Timothy C. Bell in 1991 as part of their work on data compression.

How Witten-Bell Smoothing Works?

- Witten-Bell smoothing is based on the idea of interpolating between the maximum likelihood estimate (MLE) and a fallback estimate. It assumes that the likelihood of encountering an unseen event is proportional to the number of unique events that have been observed so far.
- The probability of an unseen event is estimated using the count of unique n-grams observed in the training data. The idea is that if a particular context has produced many different n-grams, it is likely that it will also produce new, unseen ones.
- The smoothed probability for an n-gram is calculated

$$P(w_n | w_1, \dots, w_{n-1}) = \frac{C(w_1, \dots, w_n)}{C(w_1, \dots, w_{n-1}) + N(w_1, \dots, w_{n-1})} P(w_n | w_{n-1}, \dots, w_1) = \frac{C(w_1, \dots, w_{n-1})}{C(w_1, \dots, w_{n-1}) + N(w_1, \dots, w_{n-1})} C(w_1, \dots, w_n)$$

- Here:
 - $C(w_1, \dots, w_n)$ is the count of the n-gram (w_1, \dots, w_n) .
 - $C(w_1, \dots, w_{n-1})$ is the count of the (n-1)-gram prefix.
 - $N(w_1, \dots, w_{n-1})$ is the number of unique words that follow the (n-1)-gram prefix.

Witten-Bell smoothing is particularly effective for handling sparse data, which is common in NLP tasks. It smooths the probability distribution by considering both the frequency of occurrence and the diversity of possible continuations in the data.

Now, let's implement this technique using the following steps:

1. **Import Required Function:** Import the `defaultdict` function from the `collections` module.
2. **Define the Class:** Create a class called `WittenBellSmoothing` that:
 - Initializes `unigrams` and `bigrams` as default dictionaries.
 - Sets up a counter variable for the total number of unigrams.
3. **Model Training:** Train the model using a sample text corpus:
 - Split each sentence into tokens.
 - Update the unigram and bigram counts based on the tokens.
4. **Define Probability Calculation Function:** Implement a function named `bigram_prob` within the class to calculate the probability of a bigram using the Witten-Bell Smoothing technique.
5. **Create and Use Class Object:**
 - Define a sample text corpus.
 - Create an object of the `WittenBellSmoothing` class.
 - Call the `bigram_prob` function to calculate the probability for the bigram 'the cat'.

2. Jelinek-Mercer Smoothing

Jelinek-Mercer smoothing is a widely used technique in statistical language modeling, particularly in the context of n-gram models. It is a form of linear interpolation that aims to address the problem of estimating probabilities for n-grams that may not appear in the training data.

How Jelinek-Mercer Smoothing Works?

1. The core idea of Jelinek-Mercer smoothing is to combine the probability estimates from higher-order n-grams with those from lower-order n-grams (like unigram, bigram, etc.). This combination is done using a fixed weight (λ) that determines the influence of each level of n-grams.
2. The probability of an n-gram $(w_n|w_{n-1}, \dots, w_1)$ is calculated as:
$$P(w_n|w_{n-1}, \dots, w_1) = \lambda \cdot \text{PMLE}(w_n|w_{n-1}, \dots, w_1) + (1-\lambda) \cdot P(w_n|w_{n-1}, \dots, w_2)$$

$$= \lambda \cdot \text{PMLE}(w_n|w_{n-1}, \dots, w_1) + (1-\lambda) \cdot P(w_n|w_{n-1}, \dots, w_2)$$
 - Here:
 - $\text{PMLE}(w_n|w_{n-1}, \dots, w_1)$ is the maximum likelihood estimate of the n-gram.
 - λ is the smoothing parameter ($0 \leq \lambda \leq 1$) that controls the balance between the n-gram and the fallback model (lower-order n-gram).
 - $P(w_n|w_{n-1}, \dots, w_2)$ is the probability estimate from the lower-order model.
 - The process can be recursively applied to lower-order models until it reaches the unigram model.
3. **Interpretation:**
 - If $\lambda = 1$, the model fully relies on the maximum likelihood estimate of the higher-order n-gram.
 - If $\lambda = 0$, it relies entirely on the lower-order model.
 - By choosing an intermediate value for λ , the model effectively interpolates between these two extremes, allowing it to handle sparse data more robustly.
4. **Parameter Tuning:** The value of λ can be set using cross-validation or other optimization techniques to find the balance that best fits the training data. In practice, different values of λ may be used for different levels of n-grams.

Jelinek-Mercer smoothing is straightforward to implement and computationally efficient. It avoids the zero-probability problem for unseen n-grams by always falling back on lower-order estimates. The fixed λ makes it relatively simple to control the interpolation between higher and lower-order n-grams.

To implement this we will use the same process as before except we will use a lambda value in the `bigram_prob` function and calculate the probability of a given word with the overall probability using a weighted average.

Evaluating Language Models in NLP

Language modeling is the task of **predicting the next word or character** in a document and can be used to train language models that can be applied to a **wide range of natural language tasks** like text generation, text classification, and question answering.

The performance of language models in NLP is crucial to understand and can be evaluated with metrics like **perplexity, cross-entropy, and bits-per-character (BPC)**.

Introduction

Language models are very useful in a **broad range of applications** like speech recognition, machine translation part-of-speech tagging, parsing, Optical Character Recognition (OCR), handwriting recognition, information retrieval, and many other daily tasks.

- One of the main steps in the usage of language models is to **evaluate the performance beforehand** and use them in further tasks.
- This lets us build **confidence in the handling** of the language models in NLP and also lets us know if there are any places where the model may behave uncharacteristically.

In practice, we need to decide on the dataset to use, the method to evaluate, and also select a metric to evaluate language models. Let us learn about each of the elements further.

How to Evaluate a Language Model?

- **Evaluating a language model** lets us know whether one language model is better than another during experimentation and also to choose among already trained models.
- There are two ways to evaluate language models in NLP: Extrinsic evaluation and Intrinsic evaluation.
 - Intrinsic evaluation captures how well the model captures what it is supposed to capture, like probabilities.
 - Extrinsic evaluation (or task-based evaluation) captures how useful the model is in a particular task.
- **Comparing among language models:** We compare models by collecting a corpus of text which is common for models which we are comparing for.
 - We then divide the data into training and test sets and train the parameters of both models on the training set.
 - We then compare how well the two trained models fit the test set.

What Does Evaluating a Model Mean?

- After we train models, Whichever model assigns a higher probability to the test set is generally considered to accurately predicts the test set and hence a better model.

- Among multiple probabilistic language models, the better model is the one that has a tighter fit to the test data or that better predicts the details of the test data and hence will assign a higher probability to the test data.

Issue of Data Leakage or Bias in Language Models

- Most evaluation metrics for language models in NLP are based on test set probability, so it is important **not to let the test sentences into the training set**.
- Example: Assuming we are trying to compute the probability of a particular test sentence, and if our test sentence is part of the training corpus, we will **mistakenly assign** it an **artificially high probability** when it occurs in the test set.
 - We call this situation training on the test set.
 - Training on the test set introduces a bias that makes the probabilities all look too high and causes huge inaccuracies in metrics like perplexity. **

Extrinsic Evaluation

Extrinsic evaluation is the best way to evaluate the performance of a language model by **embedding** it in an application and measuring how much the application **improves**.

- It is an **end-to-end evaluation** where we can understand if a particular improvement in a component is really going to help the task at hand.
- Example: For speech recognition, we can **compare the performance of two language models** by running the speech recognizer twice, once with each language model, and seeing which gives the more accurate transcription.

Intrinsic Evaluation

We need to take advantage of intrinsic measures because **running big language models** in NLP systems end-to-end is often very **expensive**, and it is easier to have a metric that can be used to quickly evaluate potential improvements in a language model.

An intrinsic evaluation metric is one that measures the **quality of a model-independent** of any application.

- We also need a test set for an intrinsic evaluation of a language model in NLP
- The probabilities of an N-gram model training set come from the corpus it is trained on, the training set or training corpus.
- We can then measure the quality of an N-gram model by its performance on some unseen test set data called the test set or test corpus.
- We will also sometimes call test sets and other datasets that are not in our training sets held out corpora because we hold them out from the training data.

Good scores during intrinsic evaluation do not always mean better scores during extrinsic evaluation, so we need both types of evaluation in practice.

Perplexity

Perplexity is a very common method to evaluate the language model on some held-out data. It is a measure of **how well a probability model predicts a sample**.

- Perplexity is also an **intrinsic measure** (without the use of external datasets) to evaluate the performance of language models which come under NLP.
 - Perplexity as a metric quantifies **how uncertain a model is about the predictions it makes**. Low perplexity only guarantees a model is confident, not accurate.
 - Perplexity also often **correlates** well with the **model's final real-world performance**, and it can be quickly calculated using just the probability distribution the model learns from the training dataset.

The Intuition

- The basic intuition is that the **higher the perplexity measure** is, the **better** the language model is at modeling unseen sentences.
- Perplexity can also be seen as a **simple monotonic function of entropy**. But perplexity is often used instead of entropy due to the fact that it is **arguably more intuitive to our human minds than entropy**.

Calculating Perplexity

- Perplexity of a probability model like language models in NLP: For a model of an **unknown probability distribution**, and a proposed probability model, we can evaluate perplexity measure mathematically as $b^{-1/N \sum_{i=1}^N \log_b q(x_i)}$
 - We can choose b as 2.
 - In general, better models assign higher probabilities to the test events. Hence good models will have lower perplexity values and are less surprised by the test sample.
 - If all the probabilities were 1, then the perplexity would be one and the model would perfectly predict the text. Conversely, the perplexity will be higher for poorer language models.
- **Perplexity** denoted by PP of a discrete probability distribution p is mathematically defined as $PP(p) := 2^{H(p)} = 2^{-\sum_x p(x) \log_2 p(x)} = \prod_x p(x)^{-p(x)}$
 - Where H(p) is the entropy (in bits) of the distribution and x ranges over events which we will learn about further.
 - Perplexity of a random variable X may be defined as the **perplexity of the distribution over its possible values x**.
- One other formulation for Perplexity from the perspective of language models in NLP: It is the **multiplicative inverse of the probability** assigned to the test set by the language model normalized by the number of words in the test set.

-
-

We can define perplexity mathematically as:

$$PP(W) = P(w_1 w_2 \dots w_N)^{-1/N} = \frac{1}{P(w_1 w_2 \dots w_N)^{1/N}} = P(w_1 w_2 \dots w_N)^{-1/N}$$

- We know that if a language model can predict unseen words from the test set if the P(a sentence from a test set) is highest, then such a language model is more accurate.

Interpreting Perplexity

- Perplexity intuitively provides a **more human way of thinking about the random variable's uncertainty**. The reasoning is that the perplexity of a uniform discrete random variable with K outcomes is K.
 - Example: The perplexity of a fair coin is two and the perplexity of a fair six-sided die is six.
 - This kind of framework provides a frame of reference for interpreting a perplexity value.
- **Simple framework to interpret perplexity**: If the perplexity of some **random variable X is 10**, our uncertainty towards the outcome of X is equal to the **uncertainty we would feel towards a 10-sided die**, helping us intuit the uncertainty more deeply.

Perplexity to Compare Different N-Gram Models

- Steps to compute perplexity for n-gram models:
 - We first **calculate the joint probability** of all the words in the sentence under the n-gram model after we estimate the model parameters from the training corpus.
 - We will then **transform the joint probability into a perplexity for each sentence** by multiplying the probabilities of each word together.
 - We first need to calculate the length of the sentence in words by including the end-of-sentence word as well and then calculate the **perplexity = $1/(\text{pow}(\text{sentence_probability}, 1.0/\text{sentence_length}))$**
 - Then we compute a single **perplexity from the overall model** (if there are multiple sentences) as:
 - $PP(W) = 1/P(s_1 s_2 \dots s_N)$ $PP(W) = \frac{1}{P(s_1 s_2 \dots s_N)}$
- **Generic benchmarks and typical values of perplexity n-gram models**: If we assume a corpus of English with a vocabulary size of ~50,000, we can establish typical evaluation metrics for unigram, bigram, and trigram language models.
 - In a **bigram model**, **each word depends only on the previous word in the sentence** while in a **unigram model** each word is chosen completely **independently** of other words in the sentence.
 - The typical **reported perplexity figures** for such a dataset are ~74 for a trigram model, ~137 for a bigram model, and ~955 for a unigram model. The perplexity for a model that simply assigns probability 1/50,000 to each word in the vocabulary would be 50,000.
 - Hence the **trigram model gives a big improvement over bigram and unigram models** and a huge improvement over assigning a probability of 1/50,000 to each word in the vocabulary.

Perplexity in the Real World

- Perplexity is **used as a measure in training language models related to standardized datasets** like One Billion Word Benchmark. The dataset was collected from thousands of online news articles published in 2011, all broken down into their component sentences.
- Perplexity as a metric **measures how accurately a model can mimic the style of the dataset** it is being tested against models trained on datasets from some period as the benchmark dataset have an **unfair advantage due to vocabulary similarity** and may not work when testing for different time periods and slightly different datasets even though they were in the same domain.

- Perplexity **also rewards models for mimicking the test dataset**, and it may end up favoring the models most likely to imitate subtly toxic content (if the dataset is related to freely flowing language like a text from social media) as studies have shown such content is more polarized and gets easily discussed compared to non-toxic topics.

Pros and Cons

- **Advantages of using Perplexity**
 - Fast to calculate and hence allows researchers to select among models that are unlikely to perform well in real-world scenarios where computing is prohibitively costly and testing is time-consuming and expensive.
 - Useful to have an estimate of the model uncertainty/information density
- **Disadvantages of Perplexity**
 - **Not good for final evaluation** since it just measures the model's confidence and not its accuracy
 - Hard to make comparisons across different datasets with different context lengths, vocabulary sizes, word vs. character-based models, etc.
 - Perplexity can also end up rewarding models that mimic outdated datasets.

Entropy

Entropy is a metric that has been used to quantify the randomness of a process in many fields and compare worldwide languages, specifically in computational linguistics.

Definition for Entropy: The entropy (also called self-information) of a random variable is the **average level of the information**, surprise, or uncertainty inherent to the single variable's possible outcomes.

- The **more certain or the more deterministic** an event is, the **less information** it will contain. In a nutshell, the information is an increase in uncertainty or entropy.
- Entropy of a discrete distribution $p(x)$ over the event space X is given by: $H(p) = -\sum_{x \in X} p(x) \log_2 p(x)$ $H(p) = -\sum_{x \in X} p(x) \log p(x)$
 - $H(X) \geq 0$; $H(X) = 0$ only when the value of X is indeterminate and hence providing no new information
 - The smallest possible entropy for any distribution is zero.
 - We also know that the **entropy of a probability distribution is maximized when it is uniform**.

Entropy in Different Fields of NLP & AI

- In terms of probability theory NLP, language perspective, and probability theory NLP, entropy can also be defined as a statistical parameter that measures **how much information is produced for each letter of a text in the language**.
 - If the language is translated into binary digits (0 or 1) in the most efficient way, the entropy H is the average number of binary digits required per letter of the original language.
- From a **machine learning perspective**, entropy is a measure of uncertainty, and the objective of the machine learning model **is to minimize uncertainty**.
 - Decision tree learning algorithms use **relative entropy to determine the decision rules** that govern the data at each node.

- Classification algorithms in machine learning like logistic regression or artificial neural networks often employ a **standard loss function called cross entropy loss** that minimizes the average cross entropy between ground truth and predicted distributions.

Historical Perspective for Entropy

- **Entropy in Information Theory:** It was introduced by Claude Shannon in his definition is a **statistical parameter** which measures, in a certain sense, **how much information is produced** on the average for each letter of a text in the language.
 - If the language is translated into binary digits (0 or 1) in the most efficient way, the entropy is the average **number of binary digits required per letter** of the original language.
- **Entropy for a natural language:** The entropy of a natural language is the **average amount of information of one character in an infinite length of text**, which characterizes the complexity of natural language.
 - Historically, there have been many proposals for experimentally estimating the entropy rate as the **true probability distributions of natural language**.
 - Most of these approaches relied on the predictive power of humans or computational models such as n-gram language models and compression algorithms.
- **Using entropy as a metric:** The main idea is that if a model captures more of the structure of a language, then the entropy of the model should be lower and we can use entropy as a measure of the quality of the models.

Cross Entropy

Due to the fact that we **can not access an infinite amount of text** in the language, and the true distribution of the language is unknown, we **define a more useful and usable metric** called Cross Entropy.

- **Intuition for Cross entropy:** It is often used to measure the closeness of two distributions where one distribution is from the sample text (Q) that the language model aims to learn with as much proximity as possible and the other is the **empirical distribution** of the language (P).
 - Mathematical cross-entropy is defined as:
 - $H(P,Q) = E_P[-\log Q]$ $H(P,Q) = E_P[-\log Q]$ which can also be written as $H(P,Q) = H(P) + D_{KL}(P||Q)$ $H(P,Q) = H(P) + D_{KL}(P||Q)$
 - $H(P,Q)$ is the entropy and $D_{KL}(P||Q)$ is the Kullback–Leibler (KL) divergence of Q from P. It is also known as the relative entropy of P with respect to Q.
- From the formulation, we can see that the cross entropy of Q with respect to P is the **sum of two terms** entropy and relative entropy:
 - $H(P)$, the entropy of P, is the **average number of bits needed** to encode any possible outcome of P.
 - The number of **extra bits required** to encode any possible outcome of P optimized over Q.

The empirical entropy $H(P)$ is **unoptimizable**, so when we train a language model with the objective of minimizing the cross-entropy loss, the **true objective is to minimize the KL divergence** of the distribution which was learned by our language model from the empirical distribution of the language.

Handling Unknown Words

- **Tokenizers in Language Models:** Tokenization is the first and important step in any NLP pipeline, especially for language models which break unstructured data and natural language text into chunks of information that can be considered as discrete elements.
 - The **token** occurrences in a document can be used directly as a **vector** representing that document.
 - The goal when crafting the vocabulary with tokenizers is to do it in such a way that the tokenizer tokenizes as few words as possible into the unknown token.
- **Issue with unknown vocabulary/tokens:** The general approach in most tokenizers is to encode the rare words in your dataset using a special token UNK by convention so that any new out-of-vocabulary word would be labeled as belonging to the rare word category.
 - We expect the model to learn how to deal with the other words from the custom UNK token.
 - It is also **generally a bad sign** if we see that the tokenizer is producing a lot of these unknown tokens as the tokenizer was not able to retrieve a sensible representation of a word and we are losing information along the way.
- **Methods to handle unknown tokens / OOV (out of vocabulary):** Character level embeddings and sub-word tokenization are some effective ways to unknown tokens.
 - Under sub-word tokenization, **WordPiece and BPE are de facto methods** employed by successful language models such as BERT and GPT, etc.
- **Character level embeddings:** Character and subword embeddings are introduced as an attempt to **limit the size of embedding matrices** such as in BERT but they have the advantage of being able to handle new slang words, misspellings, and OOV words.
 - The required **embedding matrix is much smaller** than what is required for word-level embeddings. Generally, the vectors represent each character in any language
 - Example: Instead of a single vector for "king" like in word embeddings, there would be a separate vector for each of the letters "k", "i", "n", and "g".
 - Character embeddings do not encode the same type of information that word embeddings contain and can be thought of as encoding lexical information and may be used to enhance or enrich word-level embeddings.
 - Character level embeddings are also generally **shallow in meaning** but if we have the character embedding, every single word's vector can be formed even it is out-of-vocabulary words.
- **Subword tokenization:** Subword tokenization allows the model to have a reasonable vocabulary size while being able to learn meaningful context-independent representations and also enables the model to process words it has never seen before by **decomposing them into known subwords**.
 - Example: The word refactoring can be split into re, factor, and ing. Subwords re, factor, and ing occur more frequently than the word refactoring, and their overall meaning is also kept intact.
- **Byte-Pair Encoding (BPE):** BPE was initially developed as an algorithm to **compress texts** and then used by OpenAI for tokenization when pretraining the GPT model.
 - It is used by a lot of Transformer models like GPT, GPT-2, RoBERTa, BART, and DeBERTa.

- BPE brings the perfect balance between character and word-level hybrid representations which makes it capable of managing large corpora.
- This kind of behavior also enables the **encoding of any rare words** in the vocabulary with **appropriate subword tokens** without introducing any “unknown” tokens.

Conclusion

- Intrinsic evaluation and extrinsic evaluation are two methods to evaluate the performance of language models in NLP.
- Intrinsic evaluation captures how well the model captures what it is supposed to capture on test sets from the corpus.
- Extrinsic evaluation is also called task-based evaluation and captures how useful the model is in a particular task that is used in downstream applications.
- Entropy, Cross entropy, and Perplexity are common metrics for evaluating the performance of language models in NLP.
- Words not seen while training a language model are out of vocabulary words and can be handled using custom tokens, character level embeddings, and sub-word tokenization techniques.

Evaluation Methods in Natural Language Processing (NLP)

Evaluation in NLP (Natural Language Processing) refers to the process of assessing the quality and performance of NLP models. It involves measuring how well a model is able to complete a specific NLP task, such as text classification, sentiment analysis, machine translation, or question answering. Evaluation is typically performed using metrics that reflect the accuracy or effectiveness of the model. These metrics may vary depending on the task and the specific goals of the evaluation. For example, accuracy, precision, recall, and F1-score are common metrics for evaluating text classification and information retrieval models, while BLEU and ROUGE are metrics used in machine translation evaluation, Perplexity and WER metrics are used for Automatic speech recognition and text generation.

Let's take a deep dive in accuracy, precision, recall, F1 score and BLEU score in this part.

Accuracy, precision, recall and F1-score

These metrics are especially useful in classification tasks such as sentiment analysis, named entity recognition, and text classification. All these metrics are dependent on TP, TN, FP and FN. TP (True Positive) be the number of correctly identified positive instances, TN (True Negative) be the number of correctly identified negative instances, FP (False Positive) be the number of incorrectly identified positive instances, and FN (False Negative) be the number of incorrectly identified negative instances.

Accuracy: Accuracy is the proportion of correctly classified instances out of the total number of instances. In NLP tasks, accuracy is often used to measure the overall performance of a model.

$$Accuracy = (TP + TN) / (TP + TN + FP + FN)$$

Precision: Precision is the proportion of true positives (correctly identified instances) to the total number of instances identified as positive. In NLP tasks, precision is used to measure how many of the instances identified as positive are actually positive.

$$Precision = TP / (TP + FP)$$

Recall: Recall is the proportion of true positives to the total number of instances that are actually positive. In NLP tasks, recall is used to measure how many of the positive instances were correctly identified by the model.

$$\text{Recall} = TP / (TP + FN)$$

F1-score: F1-score is the harmonic mean of precision and recall, and is used to balance the trade-off between precision and recall. It ranges from 0 to 1, with 1 being the best possible score. In NLP tasks, F1-score is often used to evaluate the overall performance of a model.

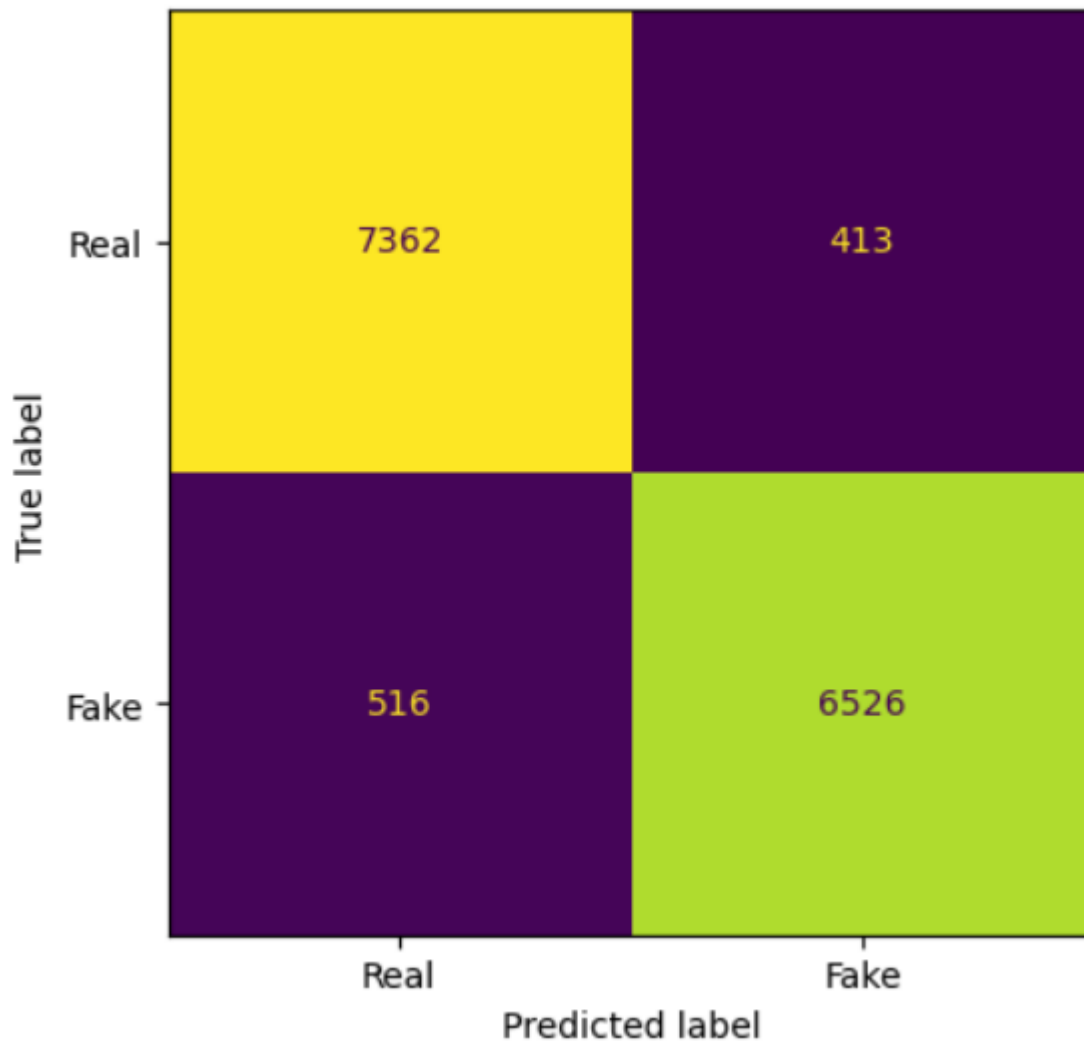
$$F1\text{-score} = 2 * (Precision * Recall) / (Precision + Recall)$$

Let's implement a text classification model and try to find out the evaluation metrics for the same. For this we are going to use the dataset from [here](#).

While running `sklearn.metrics.classification_report`, we get the following classification report:

		Per Class F1-Score			support
		precision	recall	f1-score	
Fake		0.94	0.95	0.94	7779
Real		0.94	0.93	0.94	7038
accuracy				0.94	14817
macro avg		0.94	0.94	0.94	14817
weighted avg		0.94	0.94	0.94	14817
		Average F1-Score			

We can see the recall, precision and f1-score per class and average. Instead, let us look at the **confusion matrix** for a holistic understanding of all the classes which we have taken as Y label.



The confusion matrix allows us to compute the values of True Positive (**TP**), False Positive (**FP**), and False Negative (**FN**), as shown below. Also, Let's calculate precision, recall and F1-score for each class according to formula mentioned above.

Label	TP	FP	FN	Precision (TP / (TP + FP))	Recall (TP / (TP + FN))	F1-Score (2 * (Precision * Recall) / (Precision + Recall))
Fake	7362	516	413	0.93	0.95	0.94
Real	6526	413	516	0.94	0.93	0.93

Having multiple per-class F1 scores, average of it would be better to check overall performance of a model. However, this is good when we are dealing with class imbalance or Named entity recognition task.

Let's take Macro, Weighted and Micro average of both the classes.

The macro-average is computed using the arithmetic mean of all the per-class F1 scores which treats all classes equally regardless of their support values.

Label	TP	FP	FN	Per-Class F1-Score	Macro Average
Fake	7362	516	413	0.94	$0.94 + 0.93 / 2 = 0.94$
Real	6526	413	516	0.93	

	precision	recall	f1-score	support
Fake	0.93	0.95	0.94	7775
Real	0.94	0.93	0.93	7042
accuracy			0.94	14817
macro avg	0.94	0.94	0.94	14817
weighted avg	0.94	0.94	0.94	14817

The **weighted-averaged** is calculated by taking the mean of all per-class F1 scores with consideration of each class's support.

Support mean the no. of sample used per class in calculation of classification metrics.

Label	TP	FP	FN	Per-Class F1-Score	Support	Support %	Weighted Average
Fake	7362	516	413	0.94	7775	0.52	$(0.94 \times 0.52) + (0.93 \times 0.48)$ = 0.94
Real	6526	413	516	0.93	7042	0.48	

	precision	recall	f1-score	support
Fake	0.93	0.95	0.94	7775
Real	0.94	0.93	0.93	7042
accuracy			0.94	14817
macro avg	0.94	0.94	0.94	14817
weighted avg	0.94	0.94	0.94	14817

The Micro average computes a global average by counting the sums of the True Positives (TP), False Negatives (FN), and False Positives (FP).

Label	TP	FP	FN	Micro Average
Fake	7362	516	413	$TP / (TP + ((FP+FN)/2))$
Real	6526	413	516	$13888 / (13888 + ((929+929)/2))$
Total	13888	929	929	0.94

	precision	recall	f1-score	support
Fake	0.93	0.95	0.94	7775
Real	0.94	0.93	0.93	7042
accuracy			0.94	14817
macro avg	0.94	0.94	0.94	14817
weighted avg	0.94	0.94	0.94	14817

Here, the question is we are not seeing precision and recall value for accuracy (Micro Average). Micro average is not mentioned here only the accuracy is given because micro-average essentially computes the proportion of correctly classified observations out of

all observations. In addition, if we were to do micro-average for precision and recall, we would get the same value.

Label	TP	FP	FN	Micro Average (F1- Score)	Micro Average (Precision)	Micro Average (F1- Score)
Fake	7362	516	413	$TP / (TP + ((FP+FN)/2))$	$TP / (TP + FP)$	$TP / (TP + FN)$
Real	6526	413	516	$13888 / (13888+((929+929)/2))$	$13888 / (13888+929)$	$13888 / (13888+929)$
Total	13888	929	929	0.94	0.94	0.94

Based on this we can conclude that in case of class imbalance we can deal with macro and weighted average and in case of balanced dataset overall accuracy would be fine to evaluate model performance.

[BLEU \(Bilingual Evaluation Understudy\)](#)

The BLEU score was proposed by Kishore Papineni, et al. in their 2002 paper “[BLEU: a Method for Automatic Evaluation of Machine Translation](#)”.

BLEU (Bilingual Evaluation Understudy) score is a metric used to evaluate the quality of machine-generated text or translations in Natural Language Processing (NLP). It is a statistical measure that calculates the degree of similarity between a machine-generated sentence and one or more reference sentences. The BLEU score ranges from 0 to 1, with 1 indicating perfect similarity between the machine-generated text and the reference text.

BLEU score is also used in other NLP tasks such as text summarization and image captioning. However, it is important to note that the BLEU score is not always an accurate measure of the

quality of machine-generated text, as it has certain limitations and can produce misleading results in some cases.

BLEU Score is between 0 and 1. Thus, score of >0.6 is considered the best you can achieve. Even humans would rarely achieve a perfect match. For this reason, a score closer to 1 is not realistic so it is called a model is over fitting when BLEU score is 1.

Before we get into how BLEU Score we need to understand N-grams and Precision.

N-gram

N-gram is basically a concept used in NLP for regular text processing . It is actually a set of consecutive word is a order.

For instance, in the sentence “This book is informative”, we could have n-grams such as:

- 1-gram (unigram): “This”, “book”, “is”, “informative”
- 2-gram (bigram): “This book”, “book is”, “is informative”
- 3-gram (trigram): “This book is”, “book is informative”
- 4-gram: “This book is informative”

Precision

This metric we have already discussed above. But here it works in this way.

For example, we have:

- **Target Sentence:** This book is informative
- **Predicted Sentence:** That book is informative

Precision = Number of correct predicted words / Number of total predicted words

$$Precision = 3 / 4$$

But here we need to handle it in different way. Precision is good when we would have proper class to predict. We need to handle it with the method called Clipped Precision.

Clipped Precision

Let's take an example to understand how it works.

- **Target Sentence 1:** This book is very informative
- **Predicted Sentence:** This book book is is informative

Now, we will compare each predicted word with target sentences, if words matched we will consider it as correct. We limit the count for each correct word to the maximum number of times that word occurs in the Target Sentence. Which helps to avoid word repetition.

Unigram	Predicted Count	Original Count	Clipped Count
This	1	1	1
book	2	1	1
is	2	1	1
very	0	1	0
informative	1	1	1
Total	6	4	4

Here, the word “book” occurs only once in Target Sentence. So, we have clipped the word “book” even though it occurs twice in the prediction.

Clipped Precision = Clipped number of correct predicted words / Number of total predicted words

Clipped Precision = 4/ 6

Let’s calculate BLEU score.

The mathematical representation of BLEU score is.

$$\text{BLEU} = \underbrace{\min\left(1, \exp\left(1 - \frac{\text{reference-length}}{\text{output-length}}\right)\right)}_{\text{brevity penalty}} \underbrace{\left(\prod_{i=1}^4 \text{precision}_i\right)^{1/4}}_{\text{n-gram overlap}}$$

$$\text{precision}_i = \frac{\sum_{\text{snt} \in \text{Cand-Corpus}} \sum_{i \in \text{snt}} \min(m_{\text{cand}}^i, m_{\text{ref}}^i)}{w_t^i = \sum_{\text{snt}' \in \text{Cand-Corpus}} \sum_{i' \in \text{snt}'} m_{\text{cand}}^{i'}}$$

Where,

- m_{cand}^i is the count of i-gram in candidate matching the reference translation
- m_{ref}^i is the count of i-gram in the reference translation
- w_t^i is the total number of i-grams in candidate translation

Let's take one example and calculate n-gram precision score.

Sentence : The cat is on the mat

Predicted : The cat is on the table.

As per formula mentioned above let's calculate precision

1- gram precision = $5/6$

2-gram precision = $4/5$

3-gram precision = $3/4$

4-gram precision = $2/3$

Now, we combine these scores using Geometric Average Precision (N) formula. We can use different values of N using different weights.

$$\begin{aligned}
\text{Geometric Average Precision (N)} &= \exp\left(\sum_{n=1}^N w_n \log p_n\right) \\
&= \prod_{n=1}^N p_n^{w_n} \\
&= (p_1)^{\frac{1}{4}} \cdot (p_2)^{\frac{1}{4}} \cdot (p_3)^{\frac{1}{4}} \cdot (p_4)^{\frac{1}{4}}
\end{aligned}$$

Let's dive in another step which is Brevity Penalty

Brevity Penalty

As we had seen the words like “the” and “cat” the precision is 1/1 which misleads because it encourage model to give higher score to few words. To avoid this, the Brevity Penalty penalizes sentences that are too short.

$$\text{Brevity Penalty} = \begin{cases} 1, & \text{if } c > r \\ e^{(1-r/c)}, & \text{if } c \leq r \end{cases}$$

Based on this formula is we predict very few words this Brevity Penalty is small and it can not be larger than 1 even if predicted sentence is larger than target.

Finally, BLEU Score is calculated by multiplying the Brevity Penalty with the Geometric Average of the Precision Scores.

$$\text{Bleu (N)} = \text{Brevity Penalty} \cdot \text{Geometric Average Precision Scores (N)}$$

However, there are different implementation on interent for BLEU score but the mathematics behind this is with below formula.

$$\begin{aligned}\log \text{Bleu} &= \min\left(1 - \frac{r}{c}, 0\right) + \sum_{n=1}^4 \frac{\log p_n}{4} \\ &= \min\left(1 - \frac{r}{c}, 0\right) + \frac{\log p_1 + \log p_2 + \log p_3 + \log p_4}{4}\end{aligned}$$

Finally, BLEU score is calculated mostly with N=4 where geometric average of unigram and bigram are taken into consideration.

Cons :

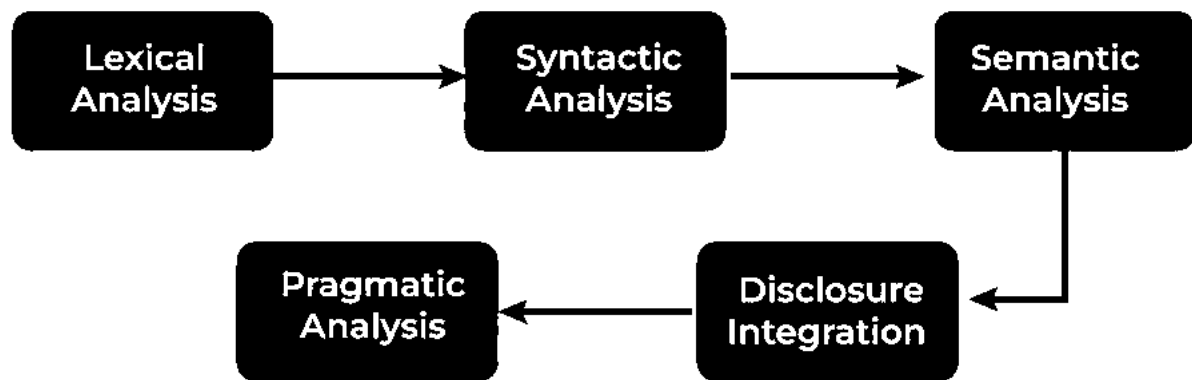
Despite its popularity, Bleu Score has been criticized for its weaknesses.

- BLEU Score does not consider the meaning of words, which can lead to incorrect assessments.
- BLEU Score only recognizes exact word matches and does not consider variants of the same word.
- BLEU Score does not distinguish between the importance of words and penalizes incorrect words equally, regardless of their relevance to the sentence.
- BLEU Score does not take into account the order of words, which can result in different sentences with the same words receiving the same (unigram) Bleu Score.

Phases of Natural Language Processing (NLP)

Natural Language Processing (NLP) is a field within artificial intelligence that allows computers to comprehend, analyze, and interact with human

language effectively. The process of NLP can be divided into five distinct phases: Lexical Analysis, Syntactic Analysis, Semantic Analysis, Discourse Integration, and Pragmatic Analysis. Each phase plays a crucial role in the overall understanding and processing of natural language.



First Phase of NLP: Lexical and Morphological Analysis

Tokenization

The lexical phase in [Natural Language Processing \(NLP\)](#) involves scanning text and breaking it down into smaller units such as paragraphs, sentences, and words. This process, known as [tokenization](#), converts raw text into manageable units called tokens or lexemes. Tokenization is essential for understanding and processing text at the word level.

In addition to tokenization, various data cleaning and feature extraction techniques are applied, including:

- **Lemmatization:** Reducing words to their base or root form.
- **Stopwords Removal:** Eliminating common words that do not carry significant meaning, such as "and," "the," and "is."
- **Correcting Misspelled Words:** Ensuring the text is free of spelling errors to maintain accuracy.

These steps enhance the comprehensibility of the text, making it easier to analyze and process.

Morphological Analysis

Morphological analysis is another critical phase in NLP, focusing on identifying morphemes, the smallest units of a word that carry meaning and cannot be further divided. Understanding morphemes is vital for grasping the structure of words and their relationships.

Types of Morphemes

1. **Free Morphemes:** Text elements that carry meaning independently and make sense on their own. For example, "bat" is a free morpheme.

2. **Bound Morphemes:** Elements that must be attached to free morphemes to convey meaning, as they cannot stand alone. For instance, the suffix "-ing" is a bound morpheme, needing to be attached to a free morpheme like "run" to form "running."

Importance of Morphological Analysis

Morphological analysis is crucial in NLP for several reasons:

- **Understanding Word Structure:** It helps in deciphering the composition of complex words.
- **Predicting Word Forms:** It aids in anticipating different forms of a word based on its morphemes.
- **Improving Accuracy:** It enhances the accuracy of tasks such as part-of-speech tagging, syntactic parsing, and machine translation.

By identifying and analyzing morphemes, the system can interpret text correctly at the most fundamental level, laying the groundwork for more advanced NLP applications.

Second Phase of NLP: Syntactic Analysis (Parsing)

Syntactic analysis, also known as parsing, is the second phase of Natural Language Processing (NLP). This phase is essential for understanding the structure of a sentence and assessing its grammatical correctness. It involves analyzing the relationships between words and ensuring their logical consistency by comparing their arrangement against standard grammatical rules.

Role of Parsing

Parsing examines the grammatical structure and relationships within a given text. It assigns [Parts-Of-Speech \(POS\) tags](#) to each word, categorizing them as nouns, verbs, adverbs, etc. This tagging is crucial for understanding how words relate to each other syntactically and helps in avoiding ambiguity. Ambiguity arises when a text can be interpreted in multiple ways due to words having various meanings. For example, the word "book" can be a noun (a physical book) or a verb (the action of booking something), depending on the sentence context.

Examples of Syntax

Consider the following sentences:

- Correct Syntax: "John eats an apple."
- Incorrect Syntax: "Apple eats John an."

Despite using the same words, only the first sentence is grammatically correct and makes sense. The correct arrangement of words according to grammatical rules is what makes the sentence meaningful.

Assigning POS Tags

During parsing, each word in the sentence is assigned a POS tag to indicate its grammatical category. Here's an example breakdown:

- Sentence: "John eats an apple."
- POS Tags:
 - John: Proper Noun (NNP)

- eats: Verb (VBZ)
- an: Determiner (DT)
- apple: Noun (NN)

Assigning POS tags correctly is crucial for understanding the sentence structure and ensuring accurate interpretation of the text.

Importance of Syntactic Analysis

By analyzing and ensuring proper syntax, NLP systems can better understand and generate human language. This analysis helps in various applications, such as machine translation, sentiment analysis, and information retrieval, by providing a clear structure and reducing ambiguity.

Third Phase of NLP: Semantic Analysis

[Semantic Analysis](#) is the third phase of Natural Language Processing (NLP), focusing on extracting the meaning from text. Unlike syntactic analysis, which deals with grammatical structure, semantic analysis is concerned with the literal and contextual meaning of words, phrases, and sentences.

Semantic analysis aims to understand the dictionary definitions of words and their usage in context. It determines whether the arrangement of words in a sentence makes logical sense. This phase helps in finding context and logic by ensuring the semantic coherence of sentences.

Key Tasks in Semantic Analysis

1. [Named Entity Recognition \(NER\)](#): NER identifies and classifies entities within the text, such as names of people, places, and organizations. These entities belong to predefined categories and are crucial for understanding the text's content.
2. [Word Sense Disambiguation \(WSD\)](#): WSD determines the correct meaning of ambiguous words based on context. For example, the word "bank" can refer to a financial institution or the side of a river. WSD uses contextual clues to assign the appropriate meaning.

Examples of Semantic Analysis

Consider the following examples:

- **Syntactically Correct but Semantically Incorrect:** "Apple eats a John."
 - This sentence is grammatically correct but does not make sense semantically. An apple cannot eat a person, highlighting the importance of semantic analysis in ensuring logical coherence.
- **Literal Interpretation:** "What time is it?"
 - This phrase is interpreted literally as someone asking for the current time, demonstrating how semantic analysis helps in understanding the intended meaning.

Importance of Semantic Analysis

Semantic analysis is essential for various NLP applications, including machine translation, information retrieval, and question answering. By ensuring that sentences are not only grammatically correct but also meaningful, semantic analysis enhances the accuracy and relevance of NLP systems.

Fourth Phase of NLP: Discourse Integration

Discourse Integration is the fourth phase of Natural Language Processing (NLP). This phase deals with comprehending the relationship between the current sentence and earlier sentences or the larger context. Discourse integration is crucial for contextualizing text and understanding the overall message conveyed.

Role of Discourse Integration

Discourse integration examines how words, phrases, and sentences relate to each other within a larger context. It assesses the impact a word or sentence has on the structure of a text and how the combination of sentences affects the overall meaning. This phase helps in understanding implicit references and the flow of information across sentences.

Importance of Contextualization

In conversations and texts, words and sentences often depend on preceding or following sentences for their meaning. Understanding the context behind these words and sentences is essential to accurately interpret their meaning.

Example of Discourse Integration

Consider the following examples:

- **Contextual Reference:** "This is unfair!"
 - To understand what "this" refers to, we need to examine the preceding or following sentences. Without context, the statement's meaning remains unclear.
- **Anaphora Resolution:** "Taylor went to the store to buy some groceries. She realized she forgot her wallet."
 - In this example, the pronoun "she" refers back to "Taylor" in the first sentence. Understanding that "Taylor" is the antecedent of "she" is crucial for grasping the sentence's meaning.

Application of Discourse Integration

Discourse integration is vital for various NLP applications, such as machine translation, sentiment analysis, and conversational agents. By understanding the relationships and context within texts, NLP systems can provide more accurate and coherent responses.

Fifth Phase of NLP: Pragmatic Analysis

Pragmatic Analysis is the fifth and final phase of Natural Language Processing (NLP), focusing on interpreting the inferred meaning of a text

beyond its literal content. Human language is often complex and layered with underlying assumptions, implications, and intentions that go beyond straightforward interpretation. This phase aims to grasp these deeper meanings in communication.

Role of Pragmatic Analysis

Pragmatic analysis goes beyond the literal meanings examined in semantic analysis, aiming to understand what the writer or speaker truly intends to convey. In natural language, words and phrases can carry different meanings depending on context, tone, and the situation in which they are used.

Importance of Understanding Intentions

In human communication, people often do not say exactly what they mean. For instance, the word "Hello" can have various interpretations depending on the tone and context in which it is spoken. It could be a simple greeting, an expression of surprise, or even a signal of anger. Thus, understanding the intended meaning behind words and sentences is crucial.

Examples of Pragmatic Analysis

Consider the following examples:

- **Contextual Greeting:** "Hello! What time is it?"
 - "Hello!" is more than just a greeting; it serves to establish contact.
 - "What time is it?" might be a straightforward request for the current time, but it could also imply concern about being late.
- **Figurative Expression:** "I'm falling for you."
 - The word "falling" literally means collapsing, but in this context, it means the speaker is expressing love for someone.

Application of Pragmatic Analysis

Pragmatic analysis is essential for applications like sentiment analysis, conversational AI, and advanced dialogue systems. By interpreting the deeper, inferred meanings of texts, NLP systems can understand human emotions, intentions, and subtleties in communication, leading to more accurate and human-like interactions.

Conclusion

The phases of NLP—Lexical Analysis, Syntactic Analysis, Semantic Analysis, Discourse Integration, and Pragmatic Analysis—each play a critical role in enabling computers to process and understand human language. By breaking down the text into manageable parts and analyzing them in different ways, NLP systems can perform complex tasks such as machine translation, sentiment analysis, and information retrieval, making significant advancements in human-computer interaction.

