## 23) Analyze the concept of N-gram models in language modeling. Discuss their strengths and weaknesses.

N-gram models are a type of statistical language model used in natural language processing (NLP) and computational linguistics. They predict the probability of a word based on the preceding n-1 words, operating under the Markov assumption, which assumes the probability of a word depends only on its immediate context of n-1 words.

An N-gram refers to a contiguous sequence of n items from a given sample of text or speech. These items can be phonemes, syllables, letters, words, or base pairs. The model's main goal is to calculate the probability of a sentence or a sequence of words.

**Types of N-grams:-**

**Ex:-The cat sat on**

- **Unigram (1-gram):** Single words; e.g., "The," "cat," "sat."

- **Bigram (2-gram):** Two consecutive words; e.g., "The cat," "cat sat."

- **Trigram (3-gram):** Three consecutive words; e.g., "The cat sat."

- **N-gram (n-gram):** N consecutive words; e.g., for n=4, "The cat sat on."

**Formula for N-gram Probability:**

$$P(w_1, w_2, \ldots, w_n) = P(w_1) \times P(w_2|w_1) \times P(w_3|w_1 w_2) \ldots P(w_n|w_{n-1})$$ Or, in gene

$$P(w_n|w_1, w_2, \ldots, w_{n-1}) \approx P(w_n|w_{n-(n-1)}, \ldots, w_{n-1})$$

### Strengths of N-gram Models:

1. **Simplicity**: N-gram models are intuitive and straightforward to implement using maximum likelihood estimation (MLE) from word counts.

2. **Efficiency**: With manageable n (e.g., 2 or 3), the models are computationally efficient, making them suitable for real-time systems like autocorrect or predictive text.

3. **Scalability for Short Contexts**: They handle short sequences effectively, especially in limited applications like phrase-level sentiment analysis.

4. **Baseline Model**: N-grams serve as a benchmark for evaluating the performance of more complex models like RNNs or Transformers.

**Example:**

For the sentence "The cat sat on the mat," a **bigram model** computes probabilities like: $P(\text{"sat"}|\text{"cat"}) = \frac{\text{Count("cat sat")}}{\text{Count("cat")}}$ If "cat sat" occurs 30 times in the corpus and "cat" occu times: $P(\text{"sat"}|\text{"cat"}) = \frac{30}{100} = 0.3$.

## Weaknesses of N-gram Models:

1. **Limited Context**: N-gram models cannot capture dependencies beyond n-1 words. For example, a **trigram model** will fail to account for relationships between non-adjacent words in sentences like:

   o "The boy who cried wolf was not believed."

2. **Data Sparsity**: Rare sequences or unseen combinations are assigned zero probability.
   **Example**: If "dog barks" never appeared in the training corpus, P("barks"|"dog") =0.

3. **Memory and Storage**: Higher-order N-grams (e.g., 4-gram or 5-gram) require substantial memory and computational resources to store probabilities for all possible word combinations.

4. **No Understanding of Syntax or Semantics**: N-grams rely purely on statistical correlations and do not understand the grammatical or semantic correctness of word sequences. For example:

   o "Mat on sat cat the" may have a higher probability than a grammatically correct sentence if trained data supports it.

5. **Scalability to Large Contexts**: As n increases, the model requires exponentially more data and computation, making it impractical for capturing long-range dependencies.

   **Example of Weakness:**

   If the sentence "The quick brown fox jumps over the lazy dog" appears once in a corpus, a trigram model would still assign P("dog"|"lazy the")=0 , showing an inability to generalize for new but plausible sequences.

### Overcoming Weaknesses:

1. **Smoothing Techniques:**

   Address the problem of zero probabilities:

   - **Laplace Smoothing**: Add 1 to all word counts to ensure no zero probabilities.

   - **Kneser-Ney Smoothing**: Improves estimates for rare n-grams.

2. **Back-off Models:**

   Use lower-order models when higher-order probabilities are zero. For example, if trigram $P(w_3|w_1, w_2)$ is unavailable, back off to bigram $P(w_3|w_2)$.

3. **Hybrid Models:**

   Combine N-gram models with neural models (e.g., embeddings) for improved performance.

4. **Alternatives:**

   Modern models like **Recurrent Neural Networks (RNNs)** or **Transformers** overcome these limitations by capturing long-range dependencies and semantic meaning.

| Strengths | Weaknesses |
|---|---|
| **Simplicity:** Easy to implement and interpret. | **Limited Context:** Cannot capture long-range dependencies beyond n-1 words. |
| **Efficiency:** Computationally efficient for small n values. | **Data Sparsity:** Assigns zero probability to unseen sequences. |
| **Baseline Model:** Useful for benchmarking advanced models. | **High Storage Requirements:** Higher-order N-grams need significant memory and storage. |
| **Works Well for Short Contexts:** Handles short, localized dependencies effectively. | **No Semantic Understanding:** Ignores grammar and meaning, relying only on statistical patterns. |
| **Scalable for Small n:** Suitable for simple tasks like autocomplete or phrase prediction. | **Scalability Issues:** Requires exponential data as n increases, making it impractical for larger contexts. |
| **Improved with Smoothing:** Techniques like Laplace smoothing mitigate some sparsity issues. | **Fails in Rare Sequences:** Still struggles with uncommon or unseen word combinations even with smoothing. |

## Example Highlight:

- **Strength:** A bigram model easily calculates $P("sat"|"cat") = \frac{\text{Count}("cat\ sat")}{\text{Count}("cat")}$.

- **Weakness:** It cannot predict relationships in "The boy who cried wolf was not believed" due to limited context.

## Applications of N-grams in NLP:

**1. Speech Recognition:** N-grams play a crucial role in modeling and recognizing spoken language patterns, improving the accuracy of speech recognition systems.

**2. Machine Translation:** In machine translation, N-grams contribute to understanding and translating phrases within a broader context, enhancing the overall translation quality.

**3. Predictive Text Input:** Predictive text input on keyboards and mobile devices relies on N-grams to suggest the next word based on the context of the input sequence.

**4. Named Entity Recognition (NER):** N-grams aid in identifying and extracting named entities from text, such as names of people, organizations, and locations.

**5. Search Engine Algorithms:** Search engines use N-grams to index and retrieve relevant documents based on user queries, improving the accuracy of search results.

## Language Model Evaluation

Language model evaluation measures how well a model predicts or generates language. It determines the model's ability to understand, predict, or generate natural language based on its task. (or) Evaluating a language model involves assessing its ability to predict a sequence of words in a language. Two common evaluation metrics are perplexity and cross-entropy .

### 1. Perplexity

- Definition:
  Perplexity quantifies how well a language model predicts a sequence of words. A lower perplexity indicates a better model since it implies higher confidence in its predictions.

- Formula:

$$\text{Perplexity} = 2^{H(p)}$$

  where $H(p)$ is the cross-entropy (average uncertainty) of the model.

- Example:
  For a test set, if the perplexity is 50, the model is as "perplexed" as guessing among 50 choices for each word. A lower perplexity (e.g., 10) means the model is better at predicting the correct word.

- Use Case:
  Comparing two models trained on the same dataset.

  - Model A: Perplexity = 25 (better).

  - Model B: Perplexity = 40.

---

### 2. Cross-Entropy

- Definition:
  Measures the divergence between the true probability distribution and the predicted distribution. Cross-entropy focuses on minimizing the uncertainty of the predictions.

- Formula:

$$H(p, q) = -\frac{1}{N} \sum_{i=1}^{N} \log(q(w_i | context))$$

  where:

  - $q(w_i | context)$: Predicted probability of word $w_i$ given its context.

  - $N$: Number of words.

- Interpretation:
  Lower cross-entropy indicates the predicted probabilities are closer to the true distribution.

- Example:
  In a sentence like "The cat sat on the mat," the actual distribution may have:

  - $P("mat") = 0.8$.
    If the model predicts $q("mat") = 0.6$, cross-entropy penalizes the deviation.

## 3. Human Evaluation

- **Definition:**

  Involves human judges assessing the quality of generated text based on fluency, coherence, and relevance.

- **Example:**

  If a model generates a summary of an article, humans might rate it based on:

  - Does it capture the main points?

  - Is the text easy to read?

- **Strength:**

  Accounts for subjective qualities like readability.

- **Weakness:**

  Expensive and time-consuming.

---

## 4. Task-Specific Evaluation

- **Definition:**

  Evaluates the model's performance on specific NLP tasks using metrics like **accuracy, precision, recall, and F1 score.**

- **Example:**

  For a sentiment analysis task:

  - Model predicts "positive" correctly for 90 out of 100 examples → Accuracy = 90%.

---

## 5. Diversity and Novelty Evaluation

- **Definition:**

  Evaluates the originality and creativity of generated text.

- **Example:**

  A chatbot generating multiple variations of responses to avoid repetition, such as:

  - "Hello, how can I help you?"

  - "Hi! What can I do for you?"

- **Use Case:**

  Important in dialogue systems to maintain engagement.

## Comparison of Methods

| Method | Strengths | Weaknesses |
|---|---|---|
| Perplexity | Provides a clear, numeric benchmark for model comparison. | Does not always correlate with task performance (e.g., text fluency). |
| Cross-Entropy | Quantifies how close predictions are to true probabilities. | Computationally expensive for large datasets. |
| Human Evaluation | Captures subjective aspects like fluency and coherence. | Time-intensive and subjective. |
| Task-Specific Metrics | Directly evaluates performance for tasks like classification or summarization. | Limited to specific tasks; may not reflect general language modeling capabilities. |
| Diversity Evaluation | Important for generating creative or varied text. | Hard to define "ideal diversity"; may penalize valid repetitions in some contexts. |

## Simplified Examples

1. **Perplexity Example:**

   A model predicts the next word in: "I am going to the..."

   - **Model A assigns probabilities:**
     $P("park") = 0.5, P("store") = 0.4, P("beach") = 0.1.$

   - **Model B assigns probabilities:**
     $P("park") = 0.3, P("store") = 0.3, P("beach") = 0.4.$

     Model A's predictions align better with actual usage, leading to **lower perplexity**.

2. **Cross-Entropy Example:**

   For "The sun is shining," actual probabilities:

   - $P("shining"|"sun is") = 0.9.$

     If the model predicts $0.7$, the penalty (negative log) reflects a mismatch.

---

**25 .Classify the process of parameter estimation in language modeling and challenges associated with estimating parameters for large language models and describe techniques?**

### Parameter Estimation in Language Modeling

**Parameter Estimation** refers to the process of finding the optimal values for the parameters of a statistical or machine learning model. In the context of language modelling , these parameters are used to compute the probabilities of word sequences.

<u>**Steps in Parameter Estimation:**</u>

1. **Define the Model**:

   o Choose the type of model, such as N-gram models, neural networks, or transformers.

   o Example: A **bigram model** requires estimating probabilities like P(w2|w1), where w1 is the preceding word.

2. **Collect Training Data**:

   o Use a large corpus of text to compute probabilities of word sequences.

   o Example: In the sentence "The cat sat on the mat," the bigram P("sat"|"cat")is estimated based on the frequency of "cat sat" in the corpus.

3. **Choose an Estimation Method**:

   o Various methods can be applied, including:

      ▪ **Maximum Likelihood Estimation (MLE)**: Finds parameters that maximize the likelihood of the observed data.

      ▪ **Bayesian Estimation**: Combines observed data with prior knowledge.

4. **Optimization**:

   o Use optimization techniques like gradient descent to minimize a loss function (e.g., cross-entropy) and adjust parameters accordingly.

**Challenges in Parameter Estimation for Large Language Models**

1. **Data Sparsity**:

   o Large vocabularies mean many word combinations are unseen during training, leading to zero probabilities.

   o **Example**: In an N-gram model, the phrase "green unicorn" might not appear in the training data, resulting in P("unicorn"|"green")=0

2. **Computational Complexity**:

   o Large models like GPT-3 have billions of parameters, requiring significant memory and computation power.

   o **Example**: Training GPT-3 on billions of tokens takes weeks on supercomputers with specialized hardware.

**3. Overfitting**: Models may memorize the training data, performing well on it but poorly on unseen data.

- **Example**: A model trained on a specific domain (e.g., medical texts) may fail on general texts.

**4**. **Optimization Challenges**: Finding the global minimum of the loss function is difficult due to high-dimensional parameter spaces.

- **Example**: Gradient descent might get stuck in local minima in complex models.

**5. Bias and Fairness**: Training data may contain biases that the model learns, affecting parameter estimates.

- **Example**: If the training data associates certain professions with a gender, the model may propagate this bias.

## Techniques for Parameter Estimation

### 1. Maximum Likelihood Estimation (MLE)

- **Definition**: Maximizes the likelihood of observing the training data given the model.
- **Example**: For a bigram model, $P(\text{"sat"}|\text{"cat"}) = \frac{\text{Count}(\text{"cat sat"})}{\text{Count}(\text{"cat"})}$.
- **Strength**: Simple and effective.
- **Weakness**: Assigns zero probabilities to unseen sequences.

---

### 2. Smoothing Techniques

- **Definition**: Adjusts probability estimates to avoid zero probabilities for unseen events.
- **Techniques**:

  1. **Laplace Smoothing (Add-One)**:

     - Adds 1 to all counts to ensure non-zero probabilities.
     - Example: If "dog barks" is unseen, $P(\text{"barks"}|\text{"dog"}) = \frac{\text{Count}(\text{"dog barks"})+1}{\text{Count}(\text{"dog"})+V}$, \
       $V$ is the vocabulary size.

  2. **Kneser-Ney Smoothing**:

     - Considers lower-order N-grams to adjust estimates for rare sequences.

### 3. Bayesian Parameter Estimation

- **Definition**: Combines observed data with prior knowledge to compute probabilities.

- **Example**: If we know the likelihood of "dog barks" from past data, Bayesian estimation adjusts it based on the new corpus.

- **Strength**: Handles uncertainty better.

- **Weakness**: Computationally expensive.

---

### 4. Regularization

- **Definition**: Penalizes overly complex models to prevent overfitting.

- **Example**: In neural language models, L2 regularization adds a penalty to the loss function for large parameter values.

---

### 5. Transfer Learning

- **Definition**: Fine-tunes a pre-trained model on a smaller, task-specific dataset.

- **Example**: GPT models are pre-trained on general corpora and fine-tuned for tasks like summarization or translation.

- **Strength**: Reduces computational cost and improves accuracy on domain-specific tasks.

**26 .Choose statistical and neural language models. Discuss how each type models language data?**

## Definition of Statistical and Neural Language Models

### Statistical Language Models

Statistical language models predict the probability of a word or sequence of words based on mathematical techniques using training data. These models rely on counting occurrences of word sequences (N-grams) and using probability distributions to predict future words.

- Example: A bigram model calculates the probability of "sat" given "cat" by:

$$P("sat"|"cat") = \frac{Count("cat\ sat")}{Count("cat")}$$

### Neural Language Models

Neural language models use artificial neural networks to learn representations of words (word embeddings) and model complex relationships in language. They capture patterns in text, such as word context, semantics, and syntactic structure.

- Example: A Recurrent Neural Network (RNN) processes sequences of words like "The dog barked" and learns how "dog" relates to "barked" using hidden layers.

# How Each Type Models Language Data

**1. Statistical Language Models**

- Working:
  Statistical models use observed frequencies of word sequences in the training data. They assume the Markov property, where the probability of a word depends only on the preceding n-1 words (N-grams).

- Steps:
  1. Count occurrences of N-grams in the corpus.
  2. Compute conditional probabilities of words given their context.
  3. Use these probabilities to predict the next word in a sequence.

- Strengths:
  - Simple and interpretable.
  - Works well with small datasets.

- Limitations:
  - Cannot capture long-range dependencies due to limited context.
  - Data sparsity leads to zero probabilities for unseen N-grams.

- Example:
  For the sentence "The cat sat on the mat," a trigram model calculates probabilities like:

$$P(\text{"mat"}|\text{"sat on the"}) = \frac{\text{Count}(\text{"sat on the mat"})}{\text{Count}(\text{"sat on the"})}$$

  If "sat on the mat" occurs 50 times and "sat on the" occurs 200 times:

$$P(\text{"mat"}|\text{"sat on the"}) = \frac{50}{200} = 0.25.$$

## 2. Neural Language Models

- Working:

  Neural language models use distributed word representations (embeddings) and neural network architectures to predict words based on the entire sequence context.

- Types:

  1. Feedforward Neural Networks: Learn relationships between words using fixed-size windows of context.

  2. Recurrent Neural Networks (RNNs): Handle sequential data and remember context across longer sequences.

  3. Transformer Models (e.g., BERT, GPT): Use self-attention mechanisms to process all words in a sentence simultaneously, capturing complex dependencies.

- Strengths:

  - Capture long-range dependencies and semantic relationships.

  - Handle larger vocabularies and complex structures.

- Limitations:

  - Computationally intensive.

  - Require large amounts of data and powerful hardware.

- Example:

  Consider the sentence "The quick brown fox jumps over the lazy dog."

  - A Transformer model like GPT assigns probabilities to each word by looking at the entire sentence simultaneously, capturing how "quick" relates to "fox" or "lazy."

## Comparison Table

| Aspect | Statistical Models | Neural Models |
|---|---|---|
| Context Captured | Short-range (e.g., n-1 words for N-grams). | Long-range context with global dependencies. |
| Handling of Data Sparsity | Requires smoothing techniques. | Handles sparsity with embeddings and pretraining. |
| Computational Complexity | Lower (simple counts). | Higher (requires GPUs/TPUs for training). |
| Performance | Baseline accuracy. | State-of-the-art performance. |
| Interpretability | Easy to interpret. | Harder to interpret due to complex networks. |

## Conclusion

- Statistical models like N-grams are effective for small-scale problems and are easy to interpret, but they struggle with complex dependencies.
- Neural models provide superior performance in capturing long-range dependencies and semantics, making them suitable for advanced tasks like translation, summarization, and chatbots.
- The choice depends on the task complexity, dataset size, and computational resources.

**27 . Categorize how these challenges affect the design and performance of language models, and explore potential solutions to these problems.**

### Challenges in Language Model Design and Performance

**1. Data Sparsity:**

- **Impact:** Many word combinations in a language are rare or unseen in the training data, leading to zero probabilities for those combinations in statistical models.

- **Solution**:

- **Smoothing Techniques:**

- Example: Add-One (Laplace) Smoothing adds a constant to all probabilities, ensuring none are zero.

- **Neural Models:**

- Use embeddings to handle rare or unseen words by representing them as continuous vectors in a high-dimensional space

2. **Long-Range Dependencies**:
- **Impact**: Statistical models (e.g., N-gram) can't capture relationships beyond a limited context (n-1 words), affecting tasks like sentiment analysis or translation where context is critical.
- **Solution**:
- **Recurrent Neural Networks (RNNs)** and **Transformers**:
- RNNs maintain context over long sequences using hidden states.
- Transformers use self-attention mechanisms to understand dependencies between all words in a sentence simultaneously.

### 3.Computational Complexity:

- **Impact:** Large language models require extensive computation and memory for training and inference, making them resource-intensive.

- **Solution:**

- **Optimization Techniques:** Reduce model size (e.g., distillation) while retaining performance.

- **Efficient Hardware:** Use GPUs or TPUs for faster computation.

### 4.Bias in Training Data:

- **Impact:** Models inherit biases (e.g., gender, cultural) present in the training data, leading to biased outputs.

- **Solution:**

- **Bias Detection and Mitigation:** Use techniques to identify and reduce biased patterns in data.

- **Curated Datasets**: Train on diverse and representative datasets.

### 5.Ethical Concerns:

- **Impact**: Risks of misuse, such as generating fake news or harmful content.

- **Solution:**

- **Monitoring and Governance:** Implement restrictions on model usage and monitor outputs.

### Computation and Memory Requirements :

**Impact :** Neural models, particularly deep learning models, require significant computational resources for training. This is especially true for models like BERT or GPT, which contain billions of parameters.

**Solution : Pretrained models** like BERT and GPT allow for transfer learning, where models are pre-trained on large corpora and then fine-tuned for specific tasks, reducing the need for massive computational resources during training.

**Model Distillation** techniques, like TinyBERT or DistilBERT , reduce the size of large models while preserving much of their performance, making them more efficient.

**Overfitting : Impact :** With large language models and a small dataset, there is a risk of the model overfitting, meaning it learns the training data too well, including noise or irrelevant patterns, which hurts its generalization ability.

**Solution** : **Regularization** methods, such as dropout and weight decay , are used to prevent overfitting.

**Cross-validation techniques** can also help ensure that the model generalizes well to unseen data.

**Flowchart Representation (if needed)**

**You can imagine the design and performance improvement process as follows:**

1.  **Data Collection:** Ensure diversity and quality of training data to minimize bias and sparsity.

2.  **Model Selection:** Choose between statistical or neural models based on task complexity.

3.  **Performance Bottlenecks**: Identify issues like data sparsity, dependency limits, or computational overhead.

4.  **Optimization and Mitigation:** Apply smoothing, embeddings, efficient architectures, or regularization.

5.  **Evaluation:** Assess fairness, accuracy, and resource efficiency.

6.  **Deployment with Safeguards:** Implement monitoring systems for ethical and responsible use.

**28 . Select the approaches used to train these models, the challenges associated with balancing performance across multiple languages.**

### Approaches to Train Language Models

Language models are trained using various techniques, depending on the model's architecture and the data availability. Some common approaches include:

1.  **Supervised Learning :-** Training on labeled datasets where input-output pairs are explicitly defined.

    **Example** : Training a named entity recognition (NER) model using a labeled dataset of sentences with entity annotations (e.g., person, location).

    **Challenges** : Requires a large volume of labeled data, which can be scarce for certain languages.

**2. Unsupervised Learning :-** Training on unlabeled data by predicting context, such as missing or next words.

-   **Example**: Models like BERT (Bidirectional Encoder Representations from Transformers) use Masked Language Modeling (MLM), where some words in a sentence are masked, and the model learns to predict them.

-   **Challenges :** While it can learn useful representations, unsupervised models often struggle with tasks requiring explicit task-specific information, such as classification or named entity recognition.

**3. Self-Supervised Learning :-** A subset of unsupervised learning where the model generates pseudo-labels from the data itself

- **Example:** GPT (Generative Pre-trained Transformer) models use causal language modeling to predict the next word based on previous words.

  **Challenges** : While highly effective, these models are often computationally intensive and require significant resources to pre-train.

  **4. Transfer Learning** :- Fine-tuning a pre-trained model on a task-specific dataset.

- **Example:** Pre-trained models like GPT-3 can be adapted for summarization, translation, or question-answering by fine-tuning.

  **Challenges** : Transfer learning can face issues when transferring knowledge between very different languages (e.g., translating from a high-resource language like English to a low-resource language like Swahili).

## 5. Semi-Supervised Learning

- **Definition**: Combines labelled and unlabelled data for training.

- **Example**: Training a model on a small labelled dataset and a larger unlabelled dataset, leveraging both resources.

## Challenges in Balancing Performance Across Multiple Languages

## 1. Vocabulary Differences

- **Challenge**: Each language has unique words, grammar, and structures.

- **Solution**:

  - Use shared embedding spaces for multilingual vocabularies.

  - Tokenize text using subword units like Byte Pair Encoding (BPE) or WordPiece, which can handle diverse languages efficiently.

## 2. Grammatical Structures

- **Challenge**: Languages differ in syntax and morphology, making it difficult to use a single model for all.

- **Solution**:

  - Incorporate language-specific layers into multilingual models to adapt to individual grammatical structures.

  - Example: Use specialized layers for syntax-heavy languages like Japanese or Korean.

### 3. Data Availability

- **Challenge**: Low-resource languages have limited training data compared to high-resource languages like English.

- **Solution**:

  - Use transfer learning to adapt models trained on high-resource languages to low-resource languages.

  - Leverage cross-lingual embeddings to share knowledge across languages.

### 4. Computational Cost

- **Challenge**: Training a multilingual model requires significant resources due to the larger dataset and complex architectures.

- **Solution**:

  - Optimize model size using techniques like pruning or knowledge distillation.

  - Use parallel processing on GPUs/TPUs to reduce training time.

### 5. Bias in Training Data

- **Challenge**: Models may favor high-resource languages, leading to poorer performance for low-resource ones.

- **Solution**:

  - Balance datasets to represent all target languages proportionally.

  - Evaluate performance equally across all languages during testing.

### Steps to Train Multilingual Models

1. **Data Preprocessing**: Tokenization (e.g., WordPiece, BPE).

   - Normalization for diverse scripts and formats.

2. **Model Training**: Use a shared architecture (e.g., mBERT or XLM-R) for multiple languages.

   - Train using tasks like next-word prediction or masked language modeling.

3. **Fine-Tuning**: Adjust the model for specific languages or tasks.

4. **Evaluation**: Test using metrics like BLEU (for translation) or perplexity across languages.

**Diagram Explanation**

**Multilingual Training Workflow**:

1. **Input**: Diverse multilingual dataset.

2. **Preprocessing**: Tokenization and normalization.

3. **Modeling**:

   - Shared embeddings.

   - Language-specific layers (optional).

4. **Training**: Multi-task learning across languages.

5. **Fine-Tuning**: Adaptation for specific languages or domains.

6. **Evaluation**: Ensure fairness and balance for all languages.