

Theoretical Foundations of Large Language Models (LLMs)

1. Language Modeling Theory

A language model assigns a probability to a sequence of words w_1, w_2, \dots, w_n .

Formally:

$$P(w_1, w_2, \dots, w_n) = P(w_t | w_1, \dots, w_{t-1})$$

This formulation is autoregressive each word is predicted based on previous context.

2. Statistical Foundations

Traditional language models include n-gram models and Markov assumptions:

$$P(w_t | w_{t-1}, \dots, w_{t-n+1})$$

These models have limitations in capturing long-range dependencies and suffer from data sparsity.

3. Neural Language Modeling

Neural networks learn to represent sequences using continuous vectors.

The objective is to minimize the cross-entropy loss:

$$L = -\log P(w_t | w_{<t})$$

4. Transformers: The Foundation of LLMs

Transformers use self-attention to process sequences in parallel.

a. Self-Attention Mechanism:

Given input X , compute queries Q , keys K , and values V :

$$Q = XW^Q, K = XW^K, V = XW^V$$

$$\text{Attention}(Q, K, V) = \text{softmax}(QK^T / \sqrt{d_k}) * V$$

b. Multi-Head Attention:

Multiple attention heads capture different relationships:

$$\text{MHA}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

5. Positional Encoding

To inject order into tokens, positional encodings are added:

$$\text{PE}(\text{pos}, 2i) = \sin(\text{pos} / 10000^{(2i/d)})$$

$$\text{PE}(\text{pos}, 2i+1) = \cos(\text{pos} / 10000^{(2i/d)})$$

6. Model Objective and Training

The model is trained to predict the next word from context:

$$L = -\log P(w_t | w_1, \dots, w_{t-1})$$

Optimized using backpropagation and gradient descent.

7. Representation Learning

Tokens are embedded into a vector space and transformed into contextual embeddings based on surrounding words.

8. Emergent Properties and Scaling Laws

LLMs improve as model size, data, and compute increase. Scaling laws suggest loss decreases as resources grow.

9. Limitations

- Non-deterministic outputs
- Training instability
- No grounding in real-world knowledge