

Algorithm: BERT (Bidirectional Encoder Representations from Transformers)

BERT is a transformer-based deep learning model developed by Google that leverages bidirectional context in text to understand language more accurately. In this project, BERT is used for multi-class classification (predicting 1–5 star ratings) based on customer review content.

How the Algorithm Works

1. Bidirectional Context Understanding

Unlike earlier models (e.g., Word2Vec, GloVe, or LSTMs), BERT reads the entire sentence in both directions (left-to-right and right-to-left). This allows it to understand the true context of each word based on surrounding words.

For example:

- In the sentence “**The bank of the river was flooded**”, BERT understands that “**bank**” refers to the **side of a river**, not a financial institution.

2. Pretraining with Two Tasks

BERT is first pretrained on massive datasets using two unsupervised tasks:

a. Masked Language Modeling (MLM)

Random words in a sentence are masked, and BERT tries to predict them.
Example:

Input: “The movie was [MASK] and entertaining.” Prediction: “great”

b. Next Sentence Prediction (NSP)

BERT learns relationships between sentences by predicting whether sentence B follows sentence A.

3. Fine-Tuning for Downstream Tasks

After pretraining, BERT is fine-tuned for specific tasks such as:

- Sentiment analysis
- Question answering
- Text classification (like review rating prediction) In fine-tuning:
 - A special token [CLS] is added at the beginning of every input.
 - BERT processes the input and outputs a contextual embedding for the [CLS] token, which represents the entire sentence.
 - This embedding is passed through a classification layer (fully connected + softmax) to predict the class (e.g., star rating).

4. How BERT Processes Review Data

1. **Tokenization:** Text is broken into subwords/tokens using BERT’s tokenizer (e.g., “unhappiness” → “un”, “##happiness”).
2. **Input Representation:**
 - [CLS] token at the start
 - [SEP] token to separate sentences (if any)
 - Positional and segment embeddings are added
3. **Transformer Encoder:** Text passes through multiple self-attention layers, which allow BERT to assign dynamic importance to each word in context.
4. **Output:** The final embedding of the [CLS] token is used for classification.

5. Final Output (In Classification)

The [CLS] token's embedding is sent to a classification head:

$$\text{Softmax}(W \cdot h_{[CLS]} + b)$$

Where:

- $h_{[CLS]}$: embedding for the [CLS] token
- W : weights of the classifier
- Output: probabilities for each rating class (0 to 4)

The class with the highest probability is selected as the prediction.

Mathematical Description of BERT

1. Token Embedding + Positional Embedding

Each token x_i in the input sequence is embedded into a dense vector e_i . BERT adds positional information:

$$h_i^0 = e_i + p_i$$

Where:

- e_i is the token embedding
- p_i is the positional embedding
- h_i is the input to the first transformer layer

2. Self-Attention Mechanism

For each input token, BERT computes **Query (Q)**, **Key (K)**, and **Value (V)** vectors:

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

Then computes scaled dot-product attention:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

Where:

- d_k is the dimension of the key vectors (used for scaling)
- The softmax ensures that attention weights sum to 1 This allows

BERT to focus on relevant words in context.

3. Multi-Head Attention

BERT uses multiple attention heads to capture different types of relationships:

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

Each head is an independent self-attention operation with different learned weights.

4. Feedforward Layer in Each Transformer Block

Each output from the attention layer goes through a fully connected feedforward network:

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

5. Final Classification Layer (Fine-Tuning Stage)

For text classification tasks, BERT uses the [CLS] token's final layer embedding $h[\text{CLS}]$:

$$\hat{y} = \text{softmax}(Wh_{[\text{CLS}]} + b)$$

Where:

- **W** and **b** are learnable parameters of the classification head
- \hat{y} is the predicted probability distribution over the classes (e.g., 5 star labels)

Key Hyperparameter:

Hyperparameter	Value
Learning rate	5e-5
Batch size	16
Epochs	6
Optimizer	AdamW
Max length	64
Loss Function	CrossEntropy

Use BERT When:

- You need **deep contextual understanding** of language (e.g., subtle sentiment or sarcasm).
- Your task requires **high accuracy** and strong generalization across diverse text.
- The text contains **ambiguous or domain-specific language** where context matters.
- You're working with a **limited labeled dataset** but want to leverage pretrained knowledge.
- You have access to **GPU or sufficient compute resources**.
- You're fine-tuning a model for a **downstream NLP task** like classification, sentiment analysis, or QA.

Avoid BERT When:

- You need **fast training or real-time inference** (e.g., high-throughput applications).
- You're working on **hardware-constrained environments** (e.g., mobile or edge devices).
- You want a **lightweight, easily interpretable model**.
- You're **rapidly prototyping** and need quick results without long training cycles.
- Your task is **simple** and doesn't require deep language understanding (e.g., keyword-based tagging).
- You're **concerned with explainability**, and need transparent feature importance (e.g., for compliance or audits).